

EDA - Exploratory Data Analysis

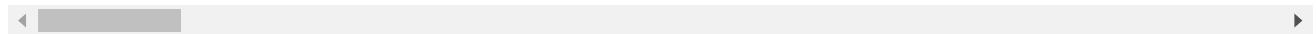
importing data

```
Loaded df_appearances from ../pickles/df_appearances.pkl
Skipping player_performance...
Skipping player_game_team_mapping...
Skipping df_games_odds...
Loaded df_teamstats from ../pickles/df_teamstats.pkl
Loaded df_shots from ../pickles/df_shots.pkl
Loaded gameresult from ../pickles/gameresult.pkl
Loaded df_after_outliers_missing from ../pickles/df_after_outliers_missing.pkl
Skipping team_performance...
Loaded df_with_categories from ../pickles/df_with_categories.pkl
Loaded df_num_after_EDA from ../pickles/df_num_after_EDA.pkl
Skipping df_games...
Loaded manipulated_data_no_outleirs from ../pickles/manipulated_data_no_outleirs.pkl
Loaded player_shots from ../pickles/player_shots.pkl
Loaded df_after_EDA from ../pickles/df_after_EDA.pkl
Loaded teamstats from ../pickles/teamstats.pkl
Loaded df_combined from ../pickles/df_combined.pkl
```

descriptive statistics

	gameID	leagueID	season	date	homeTeamID	awayTeamID	homeGoals	awayGoals	...
0	81	1	2015	2015-08-08 15:45:00	89	82	1	0	
1	82	1	2015	2015-08-08 18:00:00	73	71	0	1	
2	83	1	2015	2015-08-08 18:00:00	72	90	2	2	
3	84	1	2015	2015-08-08 18:00:00	75	77	4	2	
4	85	1	2015	2015-08-08 18:00:00	79	78	1	3	
...
12675	16131	5	2020	2021-05-23 19:00:00	168	166	1	2	
12676	16132	5	2020	2021-05-23 19:00:00	177	176	1	2	
12677	16133	5	2020	2021-05-23 19:00:00	163	235	2	0	
12678	16134	5	2020	2021-05-23 19:00:00	175	181	0	1	
12679	16135	5	2020	2021-05-23 19:00:00	225	179	1	1	

12680 rows × 47 columns



```
Defaulting to user installation because normal site-packages is not writeable
Requirement already satisfied: Autoviz in /home/leoadmin/.local/lib/python3.8/site-packages (0.1.902)
Requirement already satisfied: emoji in /home/leoadmin/.local/lib/python3.8/site-packages (from Autoviz) (2.14.1)
Requirement already satisfied: fsspec>=0.8.3 in /home/leoadmin/.local/lib/python3.8/site-packages (from Autoviz) (2025.3.0)
Requirement already satisfied: holoviews<=1.14.9 in /home/leoadmin/.local/lib/python3.8/site-packages (from Autoviz) (1.14.9)
Requirement already satisfied: hvplot~0.7.3 in /home/leoadmin/.local/lib/python3.8/site-packages (from Autoviz) (0.7.3)
Requirement already satisfied: matplotlib<=3.7.4 in /home/leoadmin/.local/lib/python3.8/site-packages (from Autoviz) (3.7.4)
Requirement already satisfied: nltk in /home/leoadmin/.local/lib/python3.8/site-packages (from Autoviz) (3.9.1)
Requirement already satisfied: numpy<1.24 in /home/leoadmin/.local/lib/python3.8/site-packages (from Autoviz) (1.23.5)
Requirement already satisfied: pandas-dq>=1.29 in /home/leoadmin/.local/lib/python3.8/site-packages (from Autoviz) (1.29)
Requirement already satisfied: pandas<2.0 in /home/leoadmin/.local/lib/python3.8/site-packages (from Autoviz) (1.5.3)
Requirement already satisfied: panel~0.14.4 in /home/leoadmin/.local/lib/python3.8/site-packages (from Autoviz) (0.14.4)
Requirement already satisfied: param==1.13.0 in /home/leoadmin/.local/lib/python3.8/site-packages (from Autoviz) (1.13.0)
Requirement already satisfied: pyamg in /home/leoadmin/.local/lib/python3.8/site-packages (from Autoviz) (5.1.0)
Requirement already satisfied: scikit-learn in /home/leoadmin/.local/lib/python3.8/site-packages (from Autoviz) (1.3.2)
Requirement already satisfied: seaborn<=0.12.2 in /home/leoadmin/.local/lib/python3.8/site-packages (from Autoviz) (0.12.2)
Requirement already satisfied: statsmodels in /home/leoadmin/.local/lib/python3.8/site-packages (from Autoviz) (0.14.1)
Requirement already satisfied: textblob in /home/leoadmin/.local/lib/python3.8/site-packages (from Autoviz) (0.18.0.post0)
Requirement already satisfied: typing-extensions>=4.1.1 in /home/leoadmin/.local/lib/python3.8/site-packages (from Autoviz) (4.12.2)
Requirement already satisfied: wordcloud in /home/leoadmin/.local/lib/python3.8/site-packages (from Autoviz) (1.9.4)
Requirement already satisfied: xgboost<1.7,>=0.82 in /home/leoadmin/.local/lib/python3.8/site-packages (from Autoviz) (1.6.2)
Requirement already satisfied: xlrd in /home/leoadmin/.local/lib/python3.8/site-packages (from Autoviz) (2.0.1)
Requirement already satisfied: pyviz-comm>=0.7.4 in /home/leoadmin/.local/lib/python3.8/site-packages (from holoviews<=1.14.9->Autoviz) (3.0.4)
Requirement already satisfied: colorcet in /home/leoadmin/.local/lib/python3.8/site-packages (from holoviews<=1.14.9->Autoviz) (3.1.0)
Requirement already satisfied: packaging in /home/leoadmin/.local/lib/python3.8/site-packages (from holoviews<=1.14.9->Autoviz) (24.2)
Requirement already satisfied: bokeh>=1.0.0 in /home/leoadmin/.local/lib/python3.8/site-packages (from hvplot~0.7.3->Autoviz) (2.4.3)
Requirement already satisfied: contourpy>=1.0.1 in /home/leoadmin/.local/lib/python3.8/site-packages (from matplotlib<=3.7.4->Autoviz) (1.1.1)
Requirement already satisfied: cycler>=0.10 in /home/leoadmin/.local/lib/python3.8/site-packages (from matplotlib<=3.7.4->Autoviz) (0.12.1)
Requirement already satisfied: fonttools>=4.22.0 in /home/leoadmin/.local/lib/python3.8/site-packages (from matplotlib<=3.7.4->Autoviz) (4.56.0)
Requirement already satisfied: kiwisolver>=1.0.1 in /home/leoadmin/.local/lib/python3.8/site-packages (from matplotlib<=3.7.4->Autoviz) (1.4.7)
```

```
Requirement already satisfied: pillow>=6.2.0 in /home/leoadmin/.local/lib/python3.8/site-packages (from matplotlib<=3.7.4->Autowiz) (10.4.0)
Requirement already satisfied: pyparsing>=2.3.1 in /home/leoadmin/.local/lib/python3.8/site-packages (from matplotlib<=3.7.4->Autowiz) (3.1.4)
Requirement already satisfied: python-dateutil>=2.7 in /home/leoadmin/.local/lib/python3.8/site-packages (from matplotlib<=3.7.4->Autowiz) (2.9.0.post0)
Requirement already satisfied: importlib-resources>=3.2.0 in /home/leoadmin/.local/lib/python3.8/site-packages (from matplotlib<=3.7.4->Autowiz) (6.4.5)
Requirement already satisfied: pytz>=2020.1 in /home/leoadmin/.local/lib/python3.8/site-packages (from pandas<2.0->Autowiz) (2025.1)
Requirement already satisfied: markdown in /home/leoadmin/.local/lib/python3.8/site-packages (from panel~=0.14.4->Autowiz) (3.7)
Requirement already satisfied: requests in /home/leoadmin/.local/lib/python3.8/site-packages (from panel~=0.14.4->Autowiz) (2.32.3)
Requirement already satisfied: tqdm>=4.48.0 in /home/leoadmin/.local/lib/python3.8/site-packages (from panel~=0.14.4->Autowiz) (4.67.1)
Requirement already satisfied: pyct>=0.4.4 in /home/leoadmin/.local/lib/python3.8/site-packages (from panel~=0.14.4->Autowiz) (0.5.0)
Requirement already satisfied: bleach in /home/leoadmin/.local/lib/python3.8/site-packages (from panel~=0.14.4->Autowiz) (6.1.0)
Requirement already satisfied: setuptools>=42 in /home/leoadmin/.local/lib/python3.8/site-packages (from panel~=0.14.4->Autowiz) (75.3.0)
Requirement already satisfied: scipy>=1.5.0 in /home/leoadmin/.local/lib/python3.8/site-packages (from scikit-learn->Autowiz) (1.10.1)
Requirement already satisfied: joblib>=1.1.1 in /home/leoadmin/.local/lib/python3.8/site-packages (from scikit-learn->Autowiz) (1.4.2)
Requirement already satisfied: threadpoolctl>=2.0.0 in /home/leoadmin/.local/lib/python3.8/site-packages (from scikit-learn->Autowiz) (3.5.0)
Requirement already satisfied: click in /usr/lib/python3/dist-packages (from nltk->Autowiz) (7.0)
Requirement already satisfied: regex>=2021.8.3 in /home/leoadmin/.local/lib/python3.8/site-packages (from nltk->Autowiz) (2024.11.6)
Requirement already satisfied: patsy>=0.5.4 in /home/leoadmin/.local/lib/python3.8/site-packages (from statsmodels->Autowiz) (1.0.1)
Requirement already satisfied: Jinja2>=2.9 in /home/leoadmin/.local/lib/python3.8/site-packages (from bokeh>=1.0.0->hvplot~=0.7.3->Autowiz) (3.1.6)
Requirement already satisfied: PyYAML>=3.10 in /usr/lib/python3/dist-packages (from bokeh>=1.0.0->hvplot~=0.7.3->Autowiz) (5.3.1)
Requirement already satisfied: tornado>=5.1 in /home/leoadmin/.local/lib/python3.8/site-packages (from bokeh>=1.0.0->hvplot~=0.7.3->Autowiz) (6.4.2)
Requirement already satisfied: zipp>=3.1.0 in /home/leoadmin/.local/lib/python3.8/site-packages (from importlib-resources>=3.2.0->matplotlib<=3.7.4->Autowiz) (3.20.2)
Requirement already satisfied: six>=1.5 in /usr/lib/python3/dist-packages (from python-dateutil>=2.7->matplotlib<=3.7.4->Autowiz) (1.14.0)
Requirement already satisfied: webencodings in /home/leoadmin/.local/lib/python3.8/site-packages (from bleach->panel~=0.14.4->Autowiz) (0.5.1)
Requirement already satisfied: importlib-metadata>=4.4 in /home/leoadmin/.local/lib/python3.8/site-packages (from markdown->panel~=0.14.4->Autowiz) (8.5.0)
Requirement already satisfied: charset-normalizer<4,>=2 in /home/leoadmin/.local/lib/python3.8/site-packages (from requests->panel~=0.14.4->Autowiz) (3.4.1)
Requirement already satisfied: idna<4,>=2.5 in /usr/lib/python3/dist-packages (from requests->panel~=0.14.4->Autowiz) (2.8)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/lib/python3/dist-packages (from requests->panel~=0.14.4->Autowiz) (1.25.8)
Requirement already satisfied: certifi>=2017.4.17 in /home/leoadmin/.local/lib/python3.8/site-packages (from requests->panel~=0.14.4->Autowiz) (2025.1.31)
Requirement already satisfied: MarkupSafe>=2.0 in /home/leoadmin/.local/lib/python3.8/site-packages (from Jinja2>=2.9->bokeh>=1.0.0->hvplot~=0.7.3->Autowiz) (2.1.5)
```

```
Shape of your Data Set loaded: (12680, 47)
#####
##### C L A S S I F Y I N G V A R I A B L E S #####
#####
Classifying variables in data set...
    Number of Numeric Columns = 16
    Number of Integer-Categorical Columns = 27
    Number of String-Categorical Columns = 1
    Number of Factor-Categorical Columns = 0
    Number of String-Boolean Columns = 0
    Number of Numeric-Boolean Columns = 0
    Number of Discrete String Columns = 0
    Number of NLP String Columns = 1
    Number of Date Time Columns = 1
    Number of ID Columns = 1
    Number of Columns to Delete = 0
47 Predictors classified...
    1 variable(s) removed since they were ID or low-information variables
    List of variables removed: ['gameID']
16 numeric variables in data exceeds limit, taking top 30 variables
    List of variables selected: ['home_xGoals', 'home_ppda', 'home_total_xAssists', 'home_total_xGoalsChain', 'home_total_xGoalsBuildup', 'home_total_blocked_shots', 'home_total_saved_shots', 'away_xGoals', 'away_ppda', 'away_total_xAssists', 'away_total_xGoalsChain', 'away_total_xGoalsBuildup', 'away_total_blocked_shots', 'home_yellowCards', 'away_yellowCards', 'away_total_saved_shots']
    Total columns > 30, too numerous to print.
To fix these data quality issues in the dataset, import FixDQ from autoviz...
    All variables classified into correct types.
```

	Data Type	Missing Values%	Unique Values%	Minimum Value	Maximum Value	DQ
gameID	int64	0.000000	100	81.000000	16135.000000	Possible ID column before modeling
leagueID	int64	0.000000	0	1.000000	5.000000	N
season	int64	0.000000	0	2014.000000	2020.000000	Possible date-time transform before modeling
date	object	0.000000	53			N
homeTeamID	int64	0.000000	1	71.000000	262.000000	Column has 55 counts greater than upper bound(256.00) or lower bound(8.00). Cap them or remove.
awayTeamID	int64	0.000000	1	71.000000	262.000000	Column has 55 counts greater than upper bound(256.00) or lower bound(8.00). Cap them or remove.
homeGoals	int64	0.000000	0	0.000000	10.000000	Column has 981 counts greater than upper bound(3.50) or lower than bound(-0.50). Cap them or remove.
awayGoals	int64	0.000000	0	0.000000	9.000000	Column has 44 counts greater than upper bound(5.00) or lower than bound(-3.00). Cap them or remove.
homeGoalsHalfTime	int64	0.000000	0	0.000000	6.000000	Column has 403 counts greater than upper bound(2.50) or lower than bound(-1.50). Cap them or remove.
awayGoalsHalfTime	int64	0.000000	0	0.000000	5.000000	Column has 225 counts greater than upper bound(2.50) or lower than bound(-1.50). Cap them or remove.
home_xGoals	float64	0.000000	NA	0.000000	6.630490	Column has 276 counts greater than upper bound(3.76) or lower than bound(-0.92). Cap them or remove.
home_shots	int64	0.000000	0	0.000000	47.000000	Column has 167 counts greater than upper bound(27.50) or lower than bound(-0.50). Cap them or remove.
home_shotsOnTarget	int64	0.000000	0	0.000000	18.000000	Column has 351 counts greater than upper bound(10.50) or lower than bound(-1.50). Cap them or remove.

	Data Type	Missing Values%	Unique Values%	Minimum Value	Maximum Value	DQ
						remove them., Column has high correlation with 'home_total_saved'. Consider dropping
home_deep	int64	0.000000	0	0.000000	42.000000	Column has 227 counts greater than upper bound(18.00) or lower than bound(-6.00). Cap them or remove
home_ppda	float64	0.000000	NA	1.897400	97.333300	Column has 540 counts greater than upper bound(21.31) or lower than bound(-1.88). Cap them or remove
home_fouls	int64	0.000000	0	0.000000	33.000000	Column has 241 counts greater than upper bound(22.50) or lower than bound(2.50). Cap them or remove
home_corners	int64	0.000000	0	0.000000	20.000000	Column has 143 counts greater than upper bound(13.00) or lower than bound(-3.00). Cap them or remove
home_yellowCards	float64	0.007886	NA	0.000000	8.000000	1 missing values. Fill them with mean, median, mode, or a constant such as 123., Column has outliers greater than upper bound (6.00) or lower bound(-2.00). Cap them or remove
home_redCards	int64	0.000000	0	0.000000	3.000000	Column has 1078 counts greater than upper bound(0.00) or lower than bound(0.00). Cap them or remove
home_total_assists	int64	0.000000	0	0.000000	8.000000	Column has 20 counts greater than upper bound(5.00) or lower than bound(-3.00). Cap them or remove them., Column has high correlation with 'homeGoals'. Consider dropping one of them
home_total_xAssists	float64	0.000000	NA	0.000000	5.512622	Column has 346 counts greater than upper bound(2.76) or lower than bound(-0.79). Cap them or remove them., Column has high correlation with 'home_xGoals'. Consider dropping one of them

	Data Type	Missing Values%	Unique Values%	Minimum Value	Maximum Value	DQ
home_total_key_passes	int64	0.000000	0	0.000000	38.000000	Column has 123 outliers greater than upper bound(22.00) or lower than lower bound(-2.00). Cap them., Column has high correlation with ['home_shots']. Consider dropping one or both.
home_total_xGoalsChain	float64	0.000000	NA	0.000000	32.011994	Column has 523 outliers greater than upper bound(10.99) or lower than lower bound(-3.97). Cap them., Column has high correlation with ['home_total_xA']. Consider dropping one or both.
home_total_xGoalsBuildup	float64	0.000000	NA	0.000000	24.437683	Column has 664 outliers greater than upper bound(6.79) or lower than lower bound(-2.88). Cap them., Column has high correlation with ['home_total_xGoalsChain']. Consider dropping one or both.
home_total_yellow_cards	int64	0.000000	0	0.000000	8.000000	Column has 20 outliers greater than upper bound(6.00) or lower than lower bound(-2.00). Cap them., Column has high correlation with ['home_yellowCards']. Consider dropping one or both.
home_total_red_cards	int64	0.000000	0	0.000000	3.000000	Column has 1064 outliers greater than upper bound(0.00) or lower than lower bound(0.00). Cap them., Column has high correlation with ['home_redCards']. Consider dropping one or both.
home_total_blocked_shots	float64	0.023659	NA	0.000000	19.000000	3 missing values. Replace them with mean, median, mode, or a constant such as 123., Column has 230 outliers greater than upper bound (9.50) or less than lower bound(0.50). Cap them or remove.
home_total_saved_shots	float64	0.023659	NA	0.000000	17.000000	3 missing values. Replace them with mean, median, mode, or a constant such as 123., Column has 532 outliers greater than upper bound (9.50) or less than lower bound(0.50). Cap them or remove.

	Data Type	Missing Values%	Unique Values%	Minimum Value	Maximum Value	DQ
						upper bound (7.00) or than lower bound(Cap them or remove
away_xGoals	float64	0.000000	NA	0.000000	6.186960	Column has 312 c greater than upper (3.10) or lower than bound(-0.91). Cap t remove
away_shots	int64	0.000000	0	0.000000	39.000000	Column has 161 c greater than upper (23.00) or lower than bound(-1.00). Cap t remove
away_shotsOnTarget	int64	0.000000	0	0.000000	15.000000	Column has 233 c greater than upper (9.50) or lower than bound(-2.50). Cap t remove them., Colum high correlati ['away_total_saved_ Consider dropping
away_deep	int64	0.000000	0	0.000000	28.000000	Column has 423 c greater than upper (13.00) or lower than bound(-3.00). Cap t remove
away_ppda	float64	0.000000	NA	2.122000	152.000000	Column has 606 c greater than upper (24.25) or lower than bound(-2.58). Cap t remove
away_fouls	int64	0.000000	0	0.000000	32.000000	Column has 81 c greater than upper (25.00) or lower than bound(1.00). Cap t remove
away_corners	int64	0.000000	0	0.000000	19.000000	Column has 283 c greater than upper (10.50) or lower than bound(-1.50). Cap t remove
away_yellowCards	float64	0.000000	NA	0.000000	9.000000	Column has 43 c greater than upper (6.00) or lower than bound(-2.00). Cap t remove
away_redCards	int64	0.000000	0	0.000000	3.000000	Column has 1396 c greater than upper (0.00) or lower than bound(0.00). Cap t remove

	Data Type	Missing Values%	Unique Values%	Minimum Value	Maximum Value	DQ
away_total_assists	int64	0.000000	0	0.000000	8.000000	Column has 790 c greater than upper (2.50) or lower than bound(-1.50). Cap t remove them., Colum high correlatio ['awayGoals']. Co dropping one o
away_total_xAssists	float64	0.000000	NA	0.000000	5.463750	Column has 389 c greater than upper (2.29) or lower than bound(-0.77). Cap t remove them., Colum high correlatio ['away_xGoals']. Co dropping one o
away_total_key_passes	int64	0.000000	0	0.000000	27.000000	Column has 169 c greater than upper (18.50) or lower than bound(-1.50). Cap t remove them., Colum high correlatio ['away_shots']. Co dropping one o
away_total_xGoalsChain	float64	0.000000	NA	0.000000	34.587459	Column has 538 c greater than upper (9.16) or lower than bound(-3.62). Cap t remove them., Colum high correlatio ['away_total_xA: Consider dropping
away_total_xGoalsBuildup	float64	0.000000	NA	0.000000	27.419105	Column has 718 c greater than upper (5.57) or lower than bound(-2.50). Cap t remove them., Colum high correlatio ['away_total_xGoalsC Consider dropping
away_total_yellow_cards	int64	0.000000	0	0.000000	9.000000	Column has 26 c greater than upper (6.00) or lower than bound(-2.00). Cap t remove them., Colum high correlatio ['away_yellowC Consider dropping
away_total_red_cards	int64	0.000000	0	0.000000	3.000000	Column has 1382 c greater than upper (0.00) or lower than bound(0.00). Cap t remove them., Colum

	Data Type	Missing Values%	Unique Values%	Minimum Value	Maximum Value	DQ
						high correlation with 'away_redCards'. Consider dropping one or both.
away_total_blocked_shots	float64	0.063091	NA	0.000000	16.000000	8 missing values. Fill them with mean, median, mode, or a constant such as 123., Column has 187 outliers greater than upper bound (8.50) or less than lower bound(-3.50). Cap them or remove.
away_total_saved_shots	float64	0.063091	NA	0.000000	13.000000	8 missing values. Fill them with mean, median, mode, or a constant such as 123., Column has 187 outliers greater than upper bound (8.50) or less than lower bound(-3.50). Cap them or remove.
gameresult	object	0.000000	0			N

Number of All Scatter Plots = 136

```
[nltk_data] Downloading collection 'popular'
[nltk_data]   | Downloading package cmudict to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package cmudict is already up-to-date!
[nltk_data]   | Downloading package gazetteers to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package gazetteers is already up-to-date!
[nltk_data]   | Downloading package genesis to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package genesis is already up-to-date!
[nltk_data]   | Downloading package gutenberg to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package gutenberg is already up-to-date!
[nltk_data]   | Downloading package inaugural to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package inaugural is already up-to-date!
[nltk_data]   | Downloading package movie_reviews to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package movie_reviews is already up-to-date!
[nltk_data]   | Downloading package names to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package names is already up-to-date!
[nltk_data]   | Downloading package shakespeare to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package shakespeare is already up-to-date!
[nltk_data]   | Downloading package stopwords to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package stopwords is already up-to-date!
[nltk_data]   | Downloading package treebank to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package treebank is already up-to-date!
[nltk_data]   | Downloading package twitter_samples to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package twitter_samples is already up-to-date!
[nltk_data]   | Downloading package omw to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package omw is already up-to-date!
[nltk_data]   | Downloading package omw-1.4 to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package omw-1.4 is already up-to-date!
[nltk_data]   | Downloading package wordnet to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package wordnet is already up-to-date!
[nltk_data]   | Downloading package wordnet2021 to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package wordnet2021 is already up-to-date!
[nltk_data]   | Downloading package wordnet31 to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package wordnet31 is already up-to-date!
[nltk_data]   | Downloading package wordnet_ic to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package wordnet_ic is already up-to-date!
[nltk_data]   | Downloading package words to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package words is already up-to-date!
[nltk_data]   | Downloading package maxent_ne_chunker to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package maxent_ne_chunker is already up-to-date!
[nltk_data]   | Downloading package punkt to
```

```
[nltk_data]      /home/leoadmin/nltk_data...
[nltk_data]      Package punkt is already up-to-date!
[nltk_data]      Downloading package snowball_data to
[nltk_data]          /home/leoadmin/nltk_data...
[nltk_data]      Package snowball_data is already up-to-date!
[nltk_data]      Downloading package averaged_perceptron_tagger to
[nltk_data]          /home/leoadmin/nltk_data...
[nltk_data]      Package averaged_perceptron_tagger is already up-
[nltk_data]          to-date!
[nltk_data]
[nltk_data] Done downloading collection popular
Could not draw wordcloud plot for date
All Plots done
Time to run AutoViz = 273 seconds
```

```
##### AUTO VISUALIZATION Completed #####
```

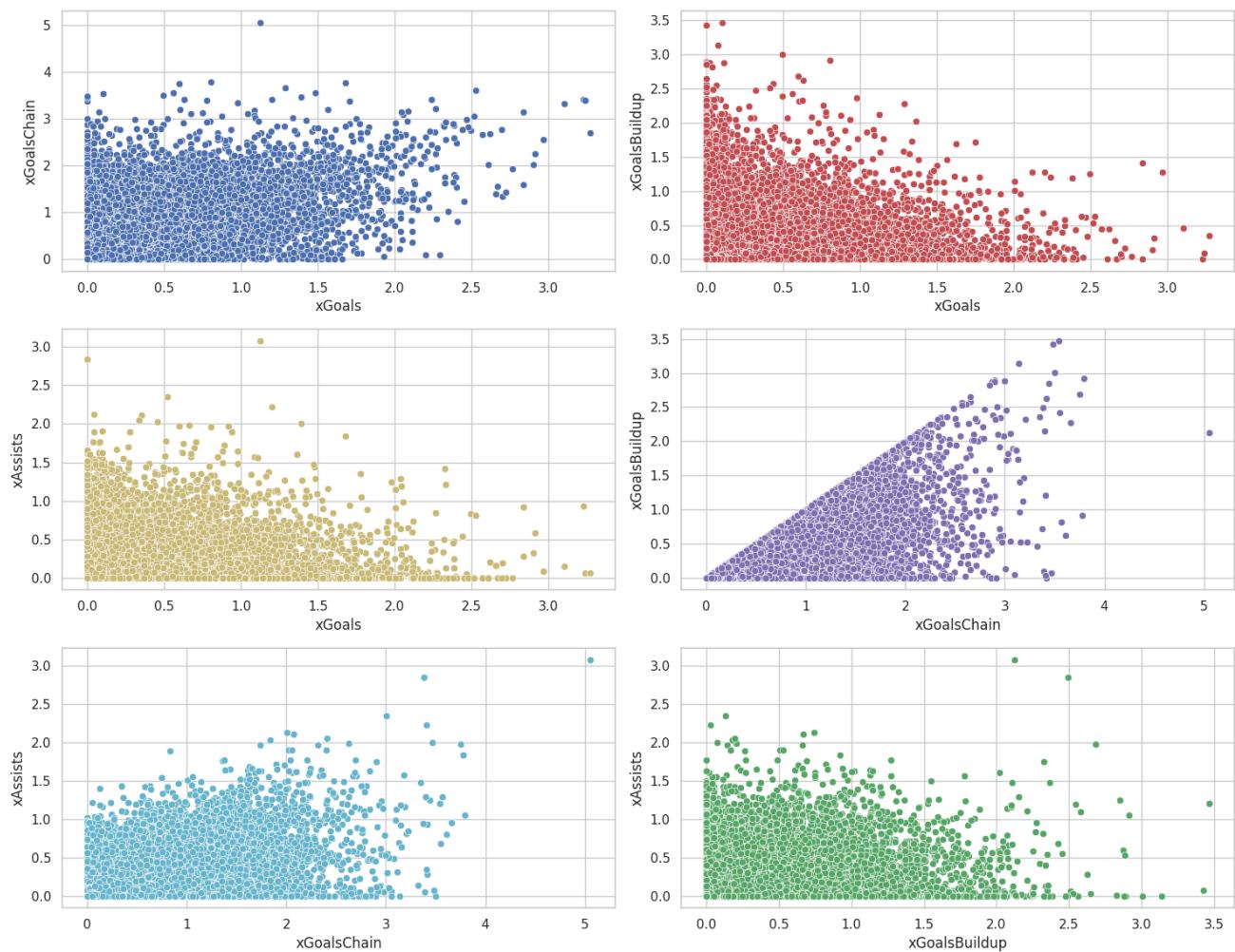
```
== Running AutoViz on df_appearances ==
max_rows_analyzed is smaller than dataset shape 356513...
    randomly sampled 150000 rows from read CSV file
Shape of your Data Set loaded: (150000, 18)
#####
##### CLASSIFYING VARIABLES #####
#####
##### Classifying variables in data set...
Number of Numeric Columns = 4
Number of Integer-Categorical Columns = 9
Number of String-Categorical Columns = 0
Number of Factor-Categorical Columns = 0
Number of String-Boolean Columns = 0
Number of Numeric-Boolean Columns = 5
Number of Discrete String Columns = 0
Number of NLP String Columns = 0
Number of Date Time Columns = 0
Number of ID Columns = 0
Number of Columns to Delete = 0
18 Predictors classified...
    No variables removed since no ID or low-information variables found in dat
a set
Since Number of Rows in data 150000 exceeds maximum, randomly sampling 150000 rows
for EDA...
To fix these data quality issues in the dataset, import FixDQ from autoviz...
    All variables classified into correct types.
```

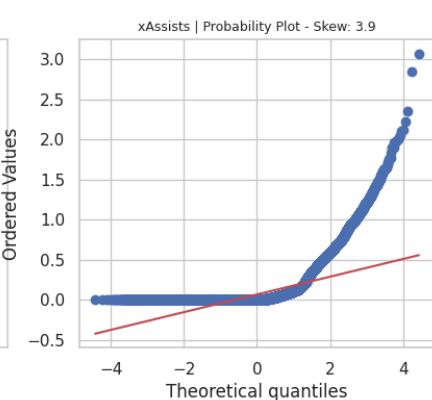
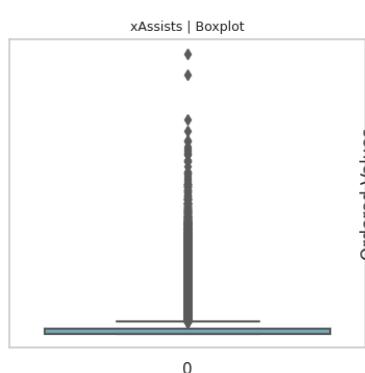
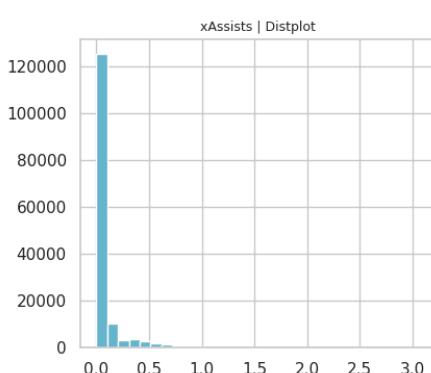
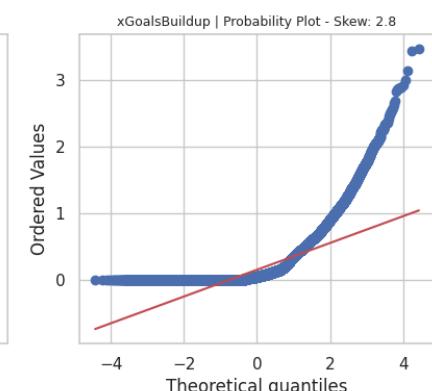
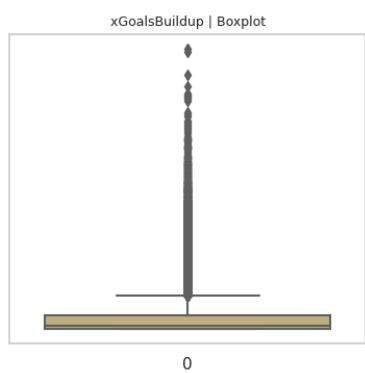
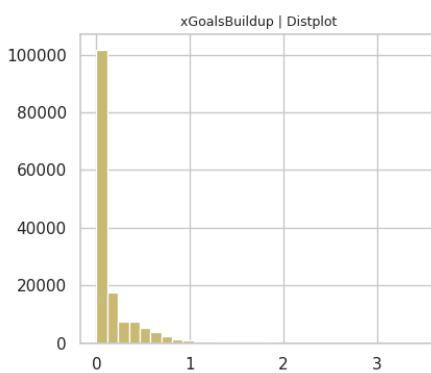
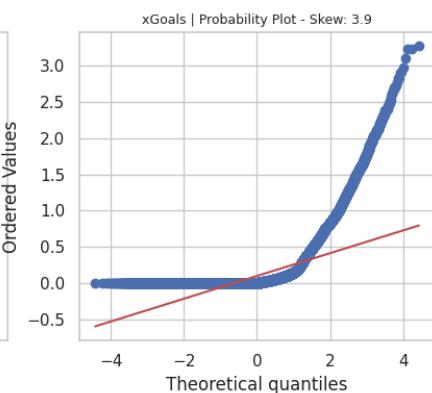
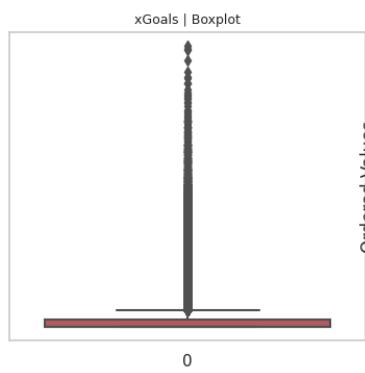
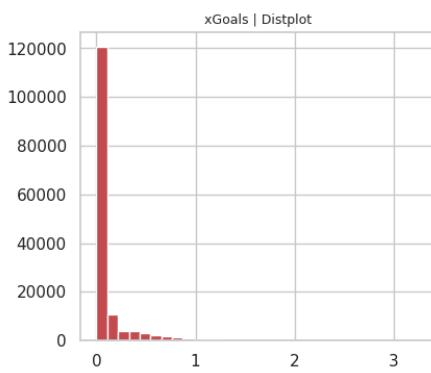
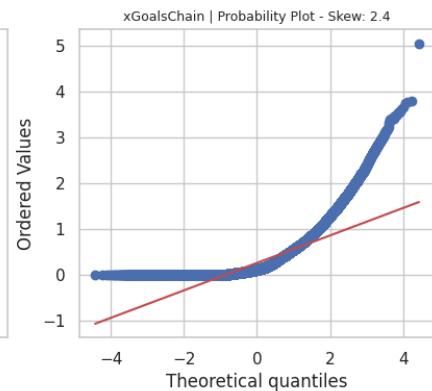
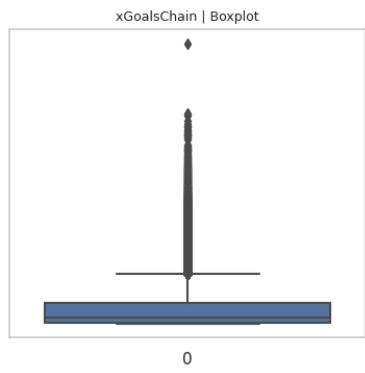
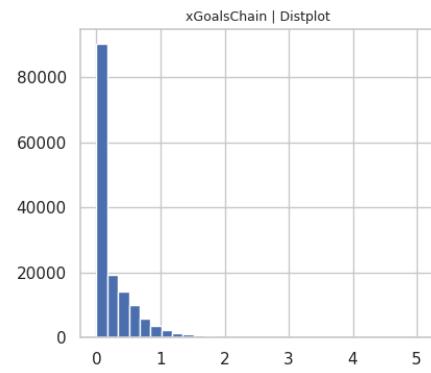
	Data Type	Missing Values%	Unique Values%	Minimum Value	Maximum Value	DQ Issue
gameID	int64	0.000000	8	81.000000	16135.000000	No issue
playerID	int64	0.000000	4	1.000000	9567.000000	Column has 3 outliers greater than upper bound (9547.50) or lower than lower bound(-4192.50). Cap them or remove them.
goals	int64	0.000000	0	0.000000	5.000000	Column has 12763 outliers greater than upper bound (0.00) or lower than lower bound(0.00). Cap them or remove them.
ownGoals	int64	0.000000	0	0.000000	2.000000	Column has 415 outliers greater than upper bound (0.00) or lower than lower bound(0.00). Cap them or remove them.
shots	int64	0.000000	0	0.000000	14.000000	Column has 15692 outliers greater than upper bound (2.50) or lower than lower bound(-1.50). Cap them or remove them.
xGoals	float64	0.000000	NA	0.000000	3.271321	Column has 19761 outliers greater than upper bound (0.19) or lower than lower bound(-0.12). Cap them or remove them.
xGoalsChain	float64	0.000000	NA	0.000000	5.052705	Column has 9034 outliers greater than upper bound (0.90) or lower than lower bound(-0.51). Cap them or remove them.
xGoalsBuildup	float64	0.000000	NA	0.000000	3.465713	Column has 18971 outliers greater than upper bound (0.41) or lower than lower bound(-0.25). Cap them or remove them.
assists	int64	0.000000	0	0.000000	4.000000	Column has 9165 outliers greater than upper bound (0.00) or lower than lower bound(0.00). Cap them or remove them.
keyPasses	int64	0.000000	0	0.000000	11.000000	Column has 9392 outliers greater than upper bound (2.50) or lower than lower bound(-1.50). Cap them or remove them.
xAssists	float64	0.000000	NA	0.000000	3.074537	Column has 18890 outliers greater than upper bound (0.14) or lower than lower bound(-0.08). Cap them or remove them.

	Data Type	Missing Values%	Unique Values%	Minimum Value	Maximum Value	DQ Issue
positionOrder	int64	0.000000	0	1.000000	17.000000	No issue
yellowCard	int64	0.000000	0	0.000000	1.000000	No issue
redCard	int64	0.000000	0	0.000000	1.000000	No issue
time	int64	0.000000	0	1.000000	90.000000	Column has 9880 outliers greater than upper bound (138.00) or lower than lower bound(10.00). Cap them or remove them.
subOut	int64	0.000000	0	0.000000	1.000000	No issue
subIn	int64	0.000000	0	0.000000	1.000000	No issue
leagueID	int64	0.000000	0	1.000000	5.000000	No issue

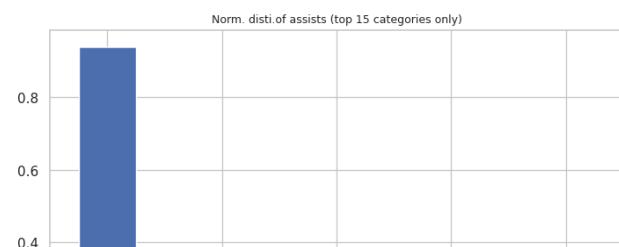
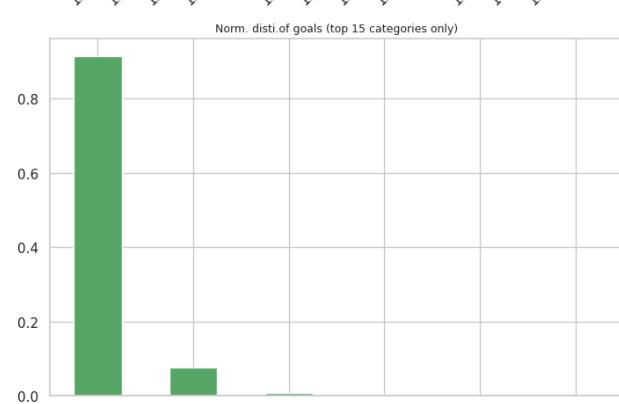
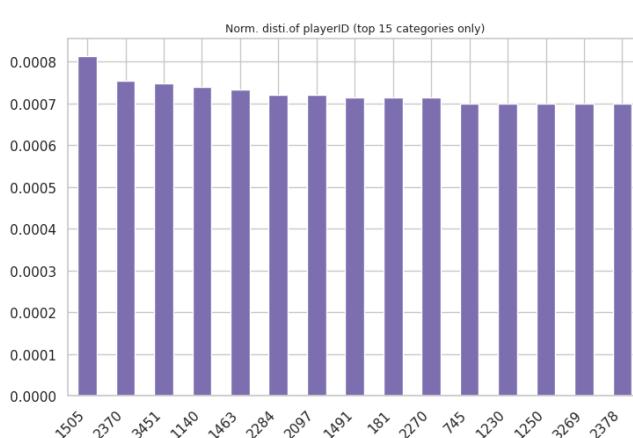
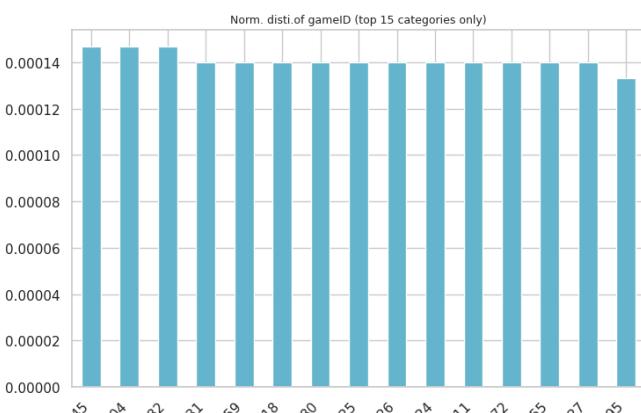
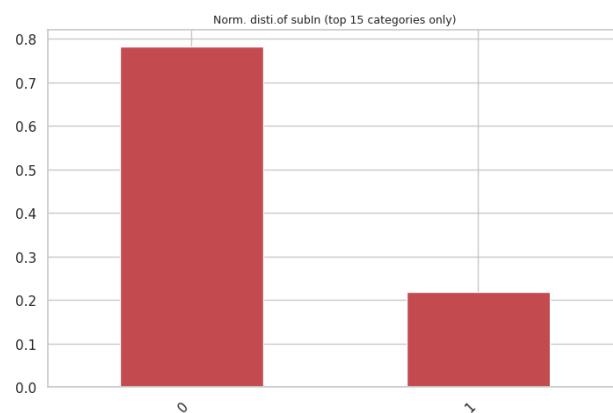
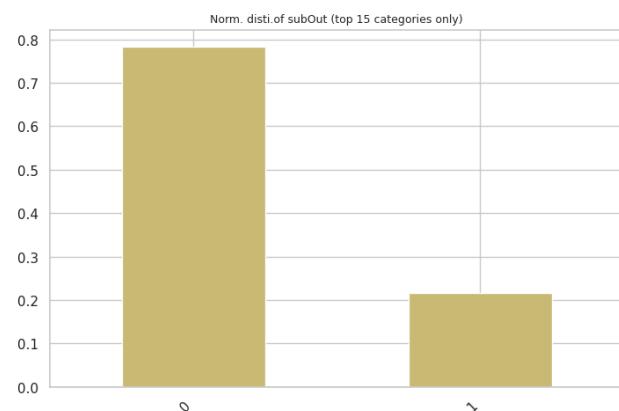
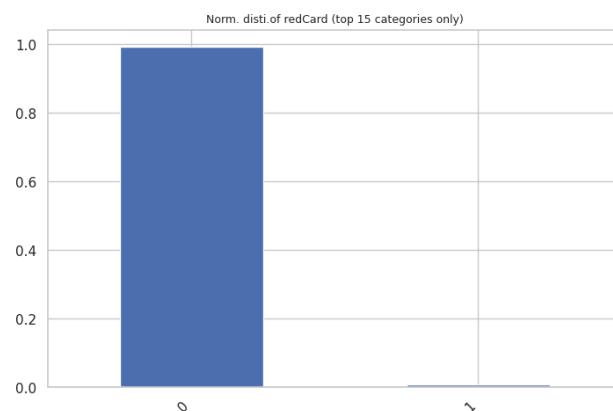
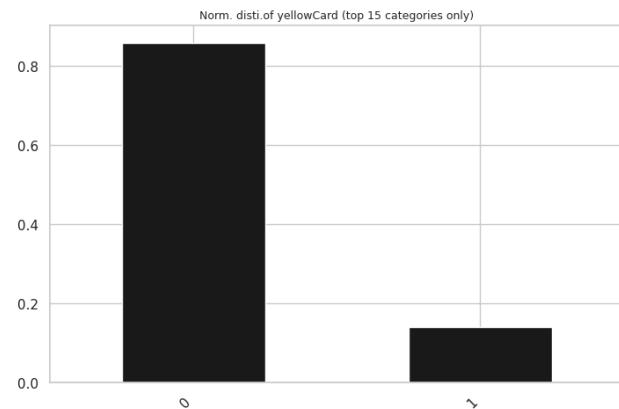
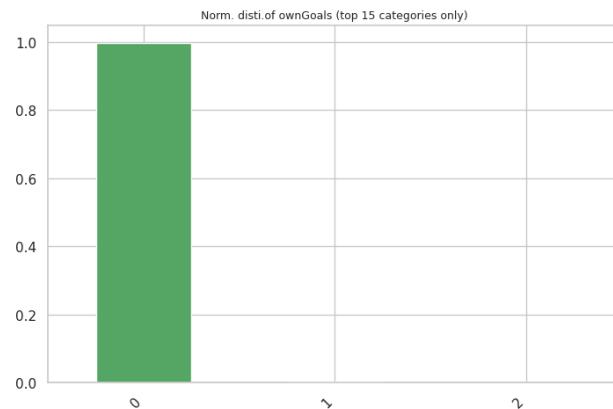
Number of All Scatter Plots = 10

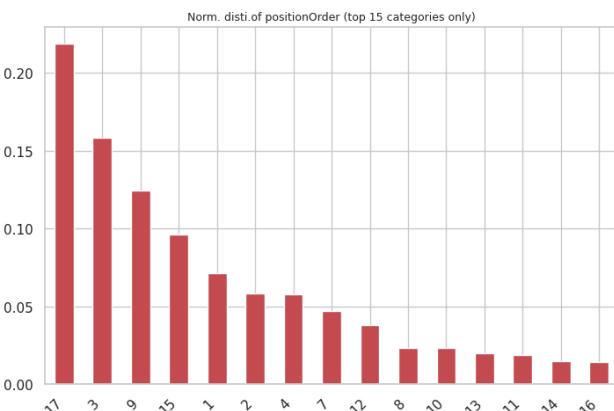
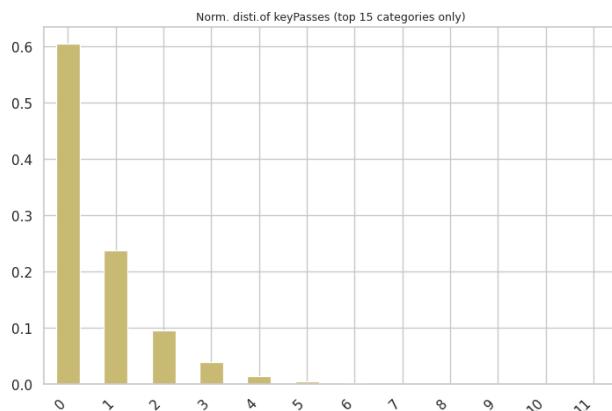
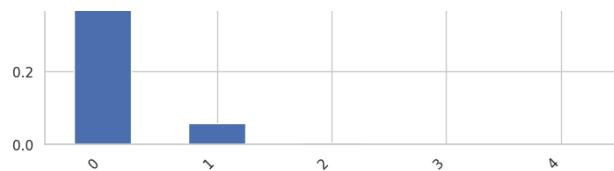
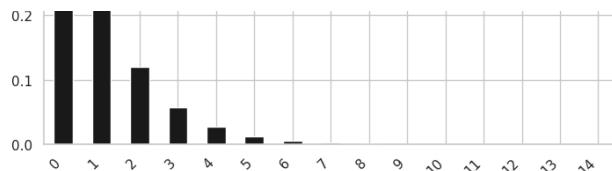
Pair-wise Scatter Plot of all Continuous Variables



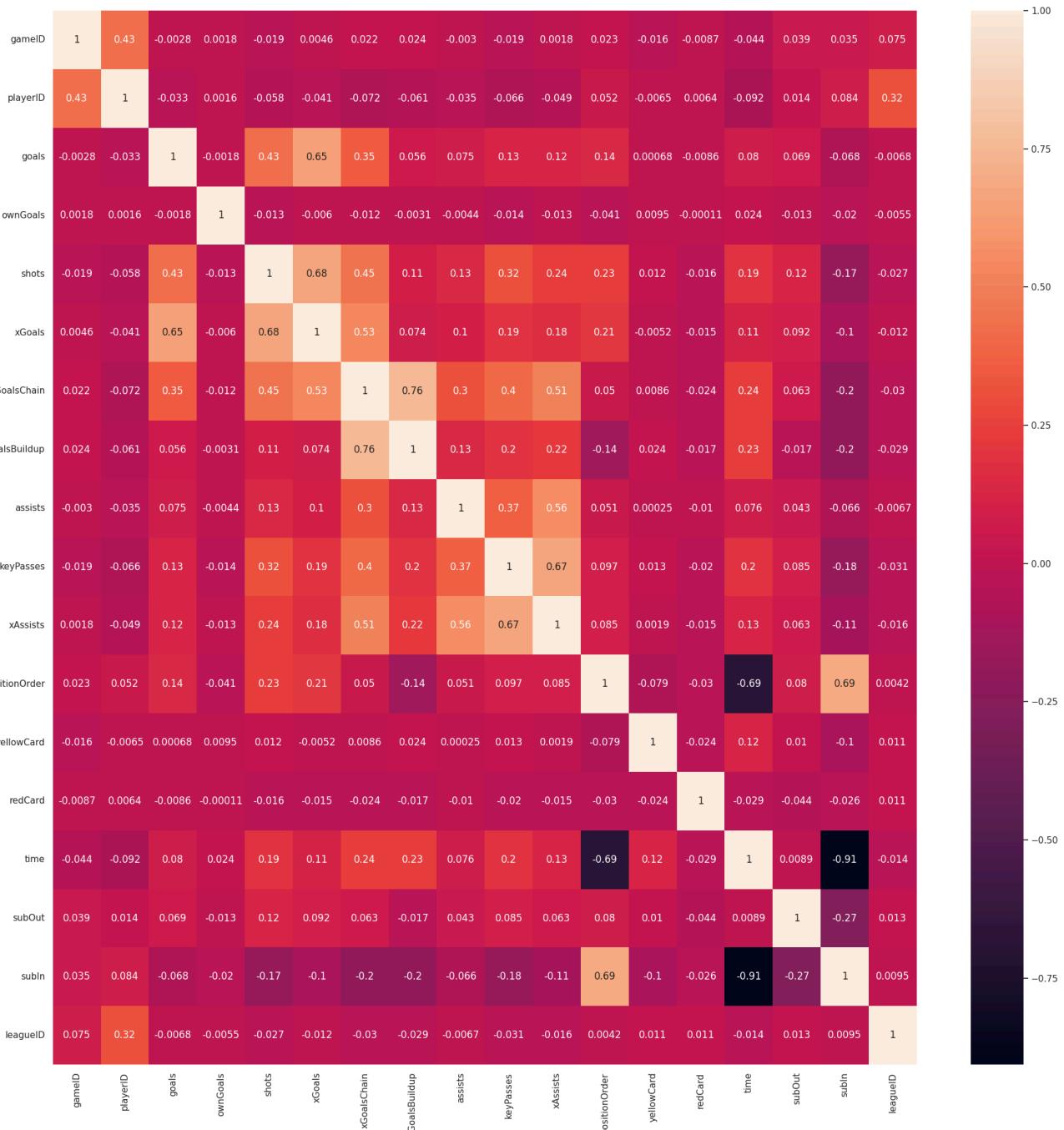


Histograms and Normalized distributions of all variables

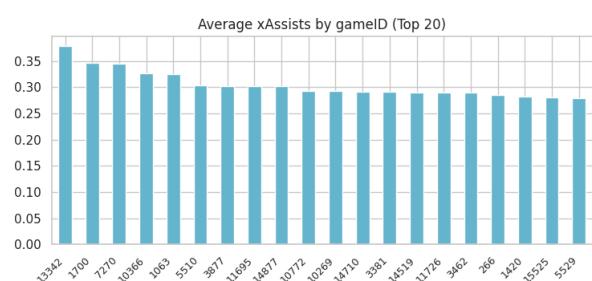
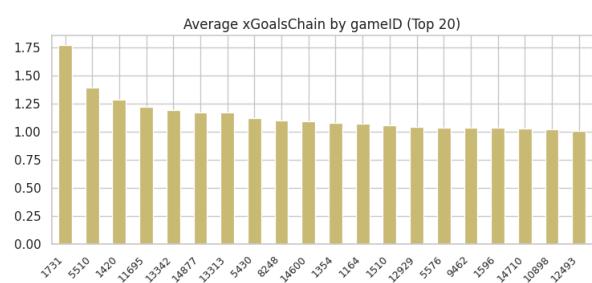
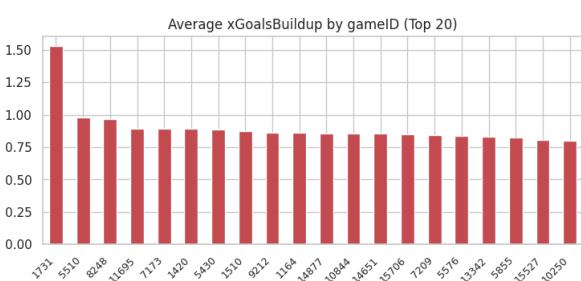
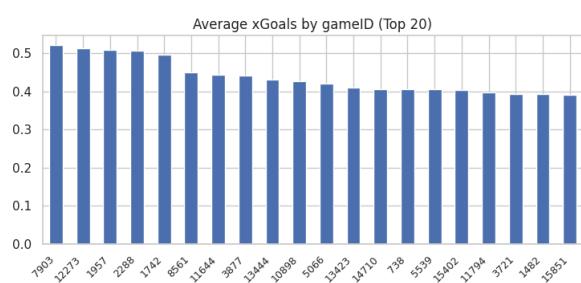


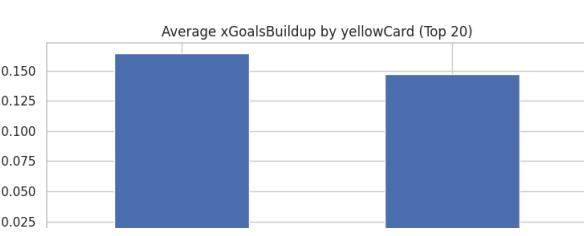
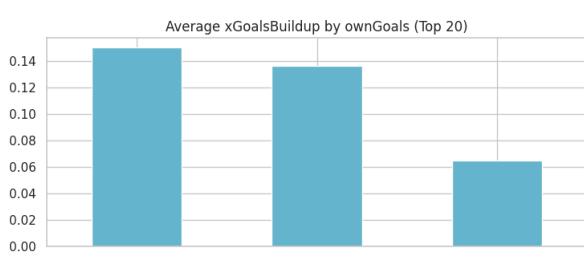
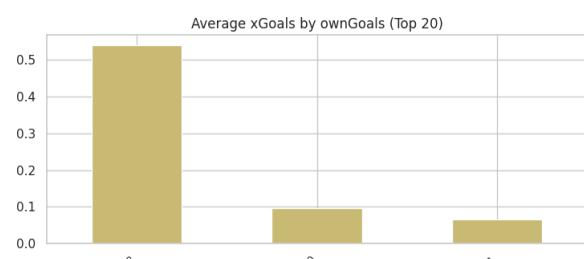
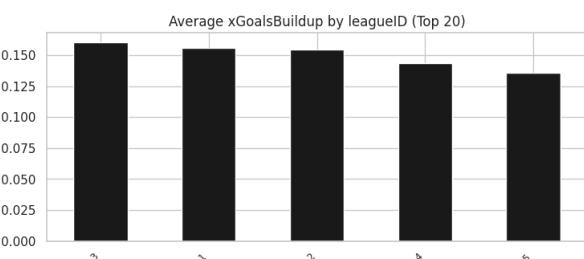
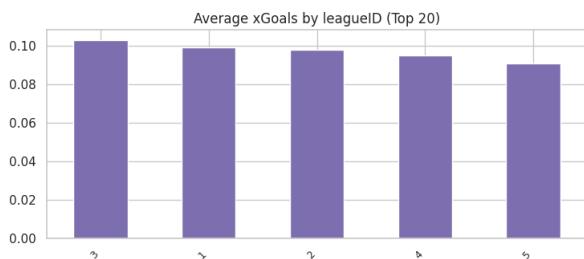
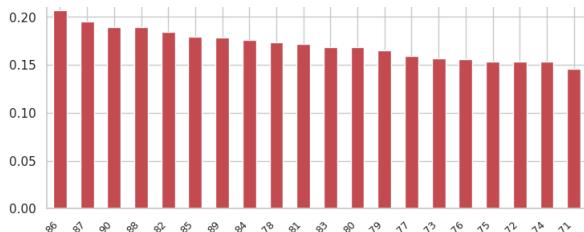
Notebook

Heatmap of all Numeric Variables including target:

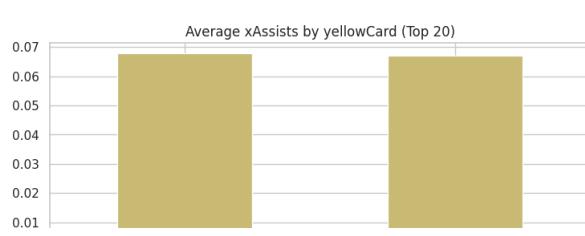
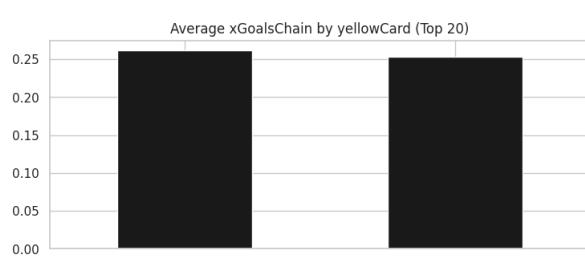
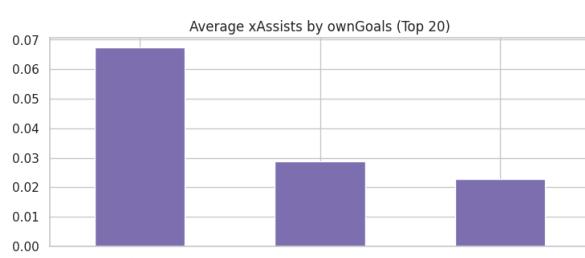
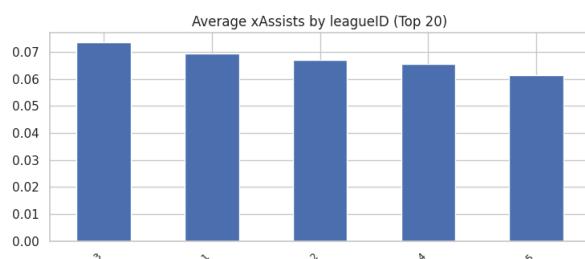
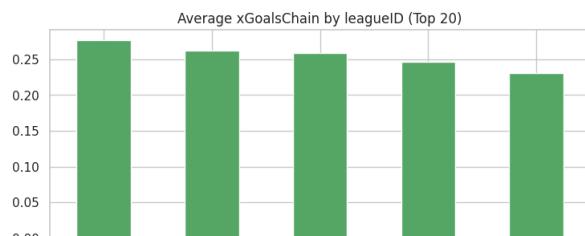
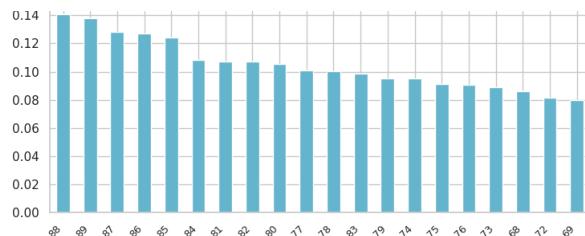


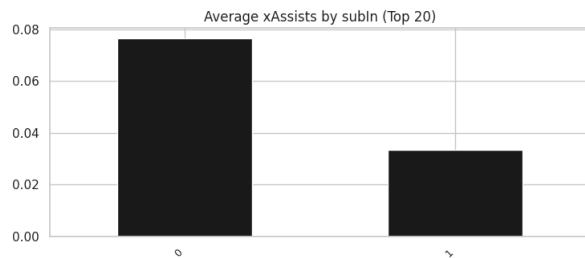
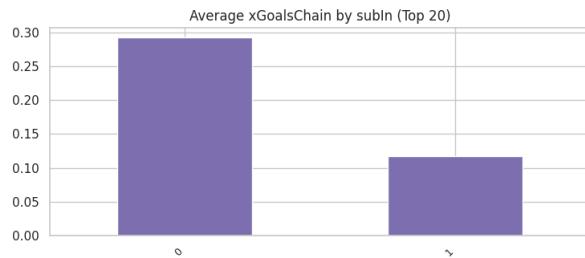
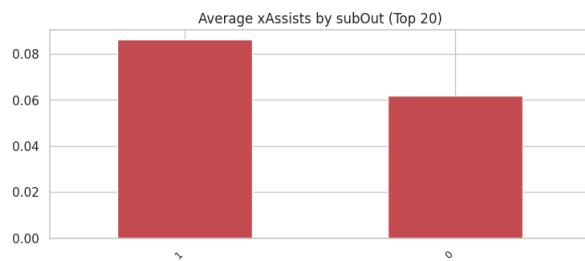
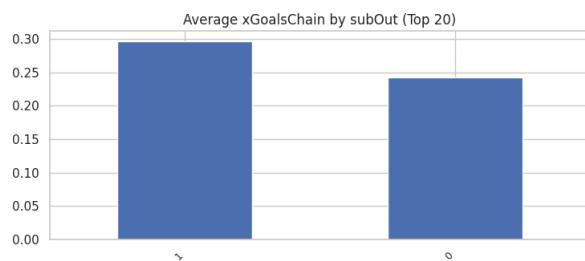
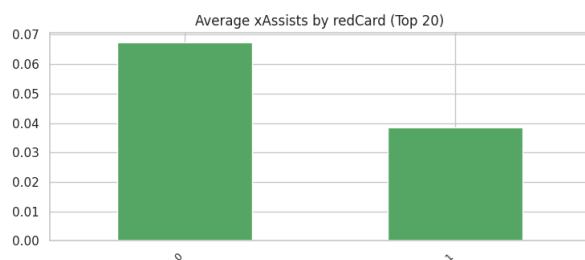
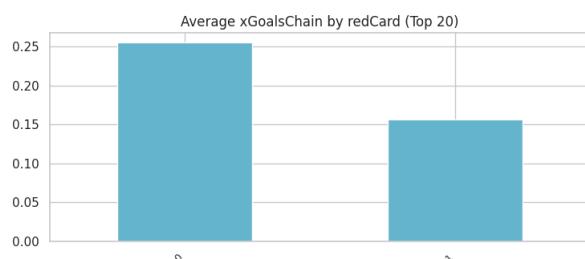
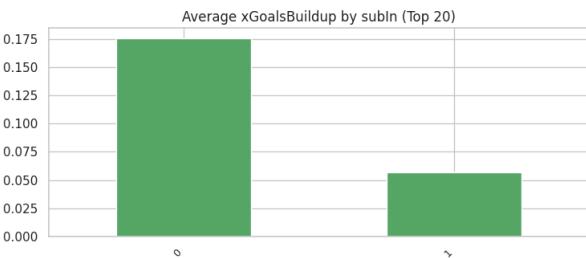
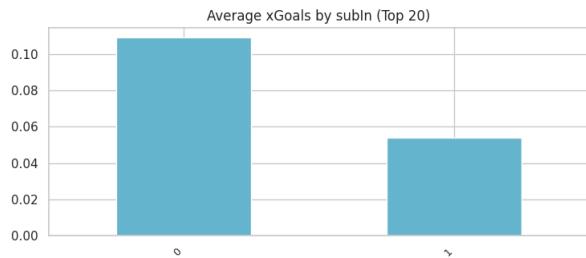
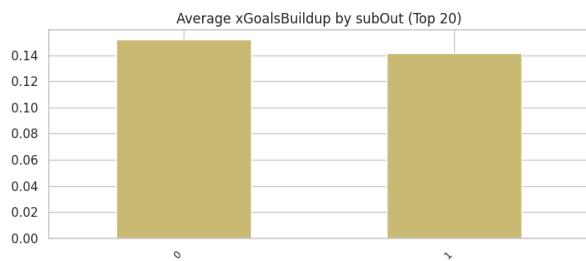
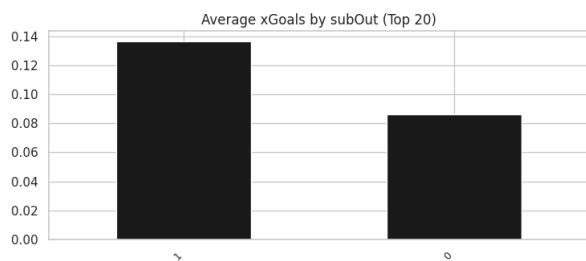
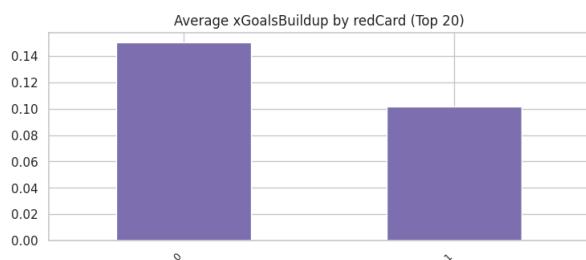
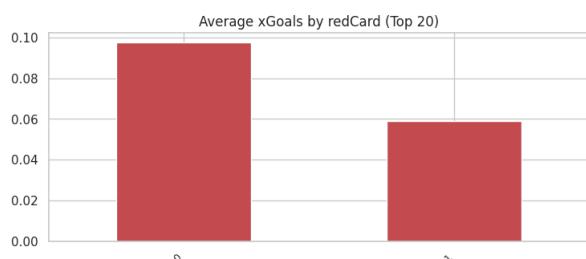
Bar plots for each Continuous by each Categorical variable





Notebook





All Plots done

Time to run AutoViz = 23 seconds

```
##### AUTO VISUALIZATION Completed #####
== Completed AutoViz on df_appearances ==
```

```
== Running AutoViz on df_teamstats ==
```

Shape of your Data Set loaded: (25360, 16)

```
#####
##### CLASSIFYING VARIABLES #####
#####
#####
```

```
#####
##### CLASSIFYING VARIABLES #####
#####
#####
```

```
#####
##### CLASSIFYING VARIABLES #####
#####
#####
```

```
#####
#####
```

Classifying variables in data set...

Number of Numeric Columns = 3

Number of Integer-Categorical Columns = 9

Number of String-Categorical Columns = 1

Number of Factor-Categorical Columns = 0

Number of String-Boolean Columns = 1

Number of Numeric-Boolean Columns = 0

Number of Discrete String Columns = 0

Number of NLP String Columns = 1

Number of Date Time Columns = 1

Number of ID Columns = 0

Number of Columns to Delete = 0

16 Predictors classified...

No variables removed since no ID or low-information variables found in data set

To fix these data quality issues in the dataset, import FixDQ from autoviz...

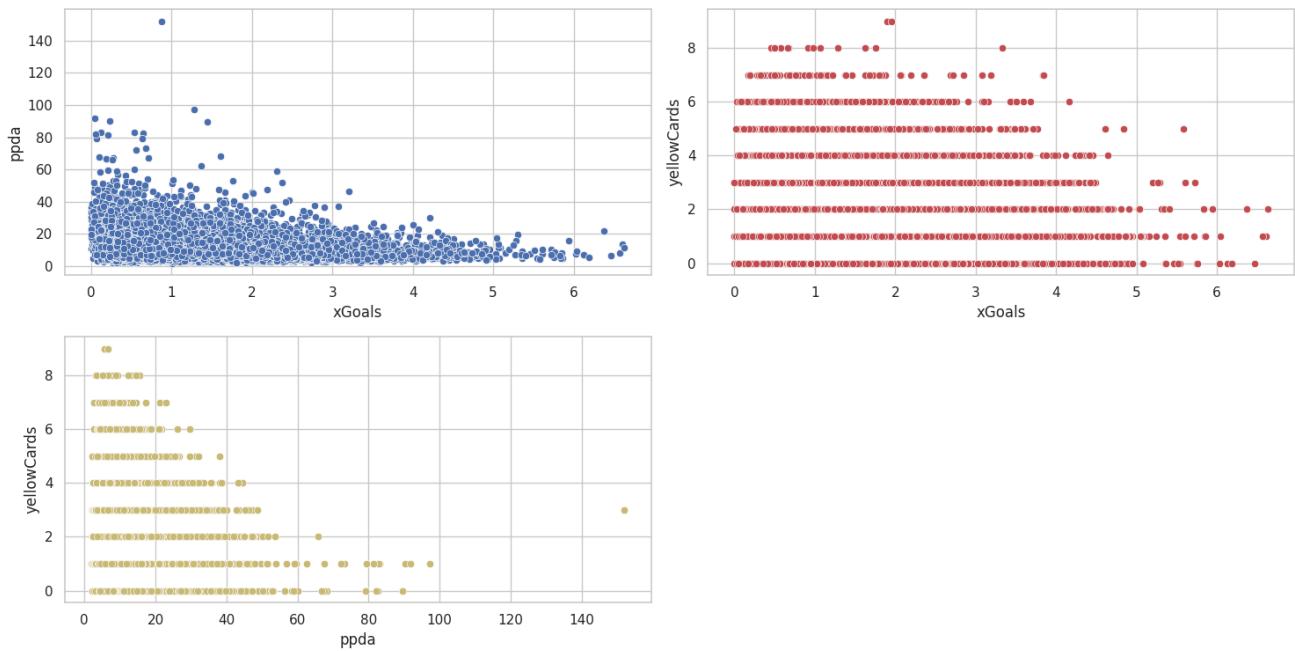
All variables classified into correct types.

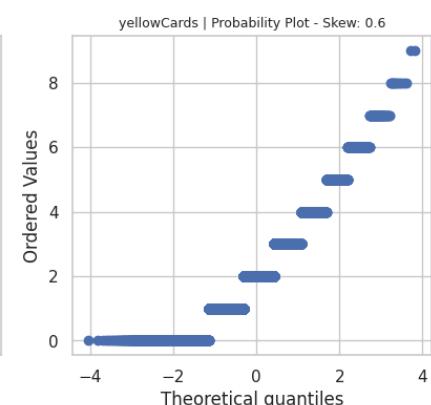
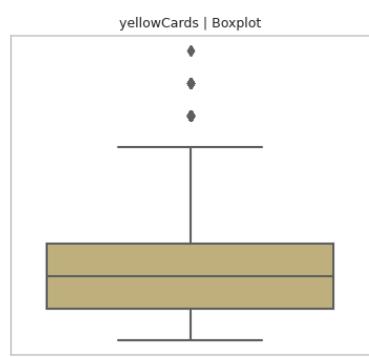
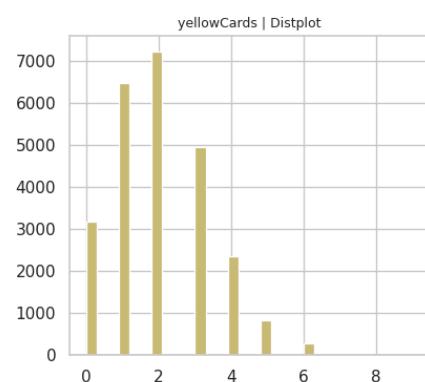
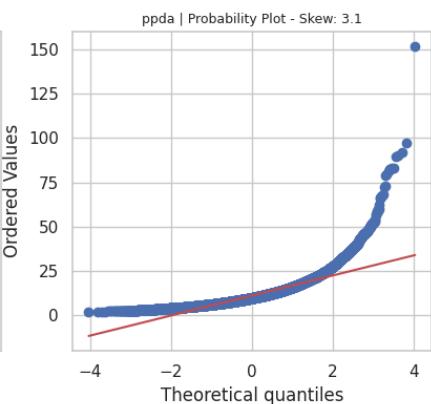
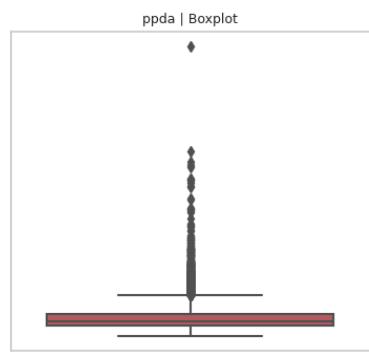
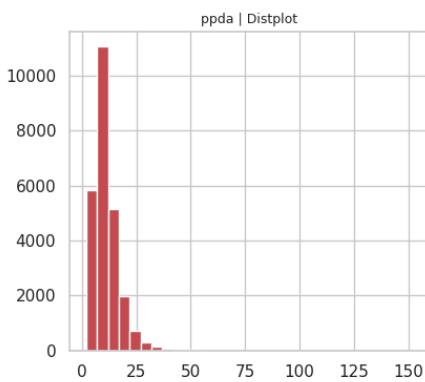
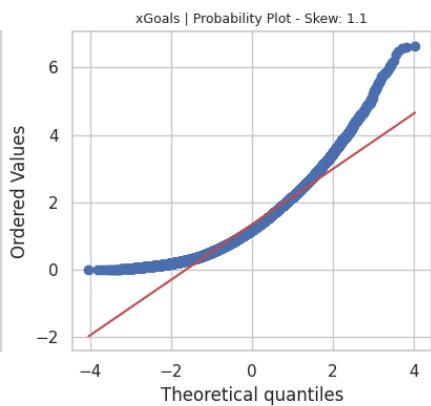
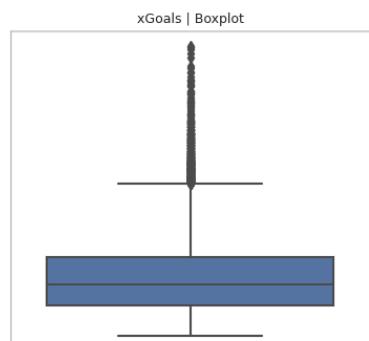
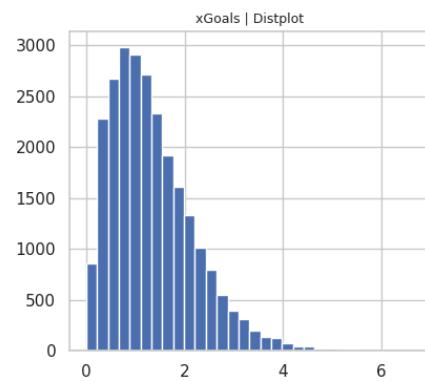
	Data Type	Missing Values%	Unique Values%	Minimum Value	Maximum Value	DQ Issue
gameID	int64	0.000000	50	81.000000	16135.000000	No issue
teamID	int64	0.000000	0	71.000000	262.000000	Column has 110 outliers greater than upper bound (256.00) or lower than lower bound(8.00). Cap them or remove them.
season	int64	0.000000	0	2014.000000	2020.000000	Possible date-time colum: transform before modeling step.
date	object	0.000000	26			No issue
location	object	0.000000	0			No issue
goals	int64	0.000000	0	0.000000	10.000000	Column has 156 outliers greater than upper bound (5.00) or lower than lower bound(-3.00). Cap them or remove them.
xGoals	float64	0.000000	NA	0.000000	6.630490	Column has 556 outliers greater than upper bound (3.48) or lower than lower bound(-0.97). Cap them or remove them.
shots	int64	0.000000	0	0.000000	47.000000	Column has 263 outliers greater than upper bound (26.50) or lower than lower bound(-1.50). Cap them or remove them.
shotsOnTarget	int64	0.000000	0	0.000000	18.000000	Column has 469 outliers greater than upper bound (10.50) or lower than lower bound(-1.50). Cap them or remove them.
deep	int64	0.000000	0	0.000000	42.000000	Column has 722 outliers greater than upper bound (15.50) or lower than lower bound(-4.50). Cap them or remove them.
ppda	float64	0.000000	NA	1.897400	152.000000	Column has 1177 outliers greater than upper bound (22.85) or lower than lower bound(-2.34). Cap them or remove them.
fouls	int64	0.000000	0	0.000000	33.000000	Column has 139 outliers greater than upper bound (25.00) or lower than lower bound(1.00). Cap them or remove them.
corners	int64	0.000000	0	0.000000	20.000000	Column has 178 outliers greater than upper bound (13.00) or lower than lower bound(-3.00). Cap them or remove them.

	Data Type	Missing Values%	Unique Values%	Minimum Value	Maximum Value	DQ Issue
yellowCards	float64	0.003943	NA	0.000000	9.000000	1 missing values. Impute them with mean, median, mode, or a constant value such as 123., Column has 80 outliers greater than upper bound (6.00) or lower than lower bound(-2.00). Cap them or remove them.
redCards	int64	0.000000	0	0.000000	3.000000	Column has 2474 outliers greater than upper bound (0.00) or lower than lower bound(0.00). Cap them or remove them.
result	object	0.000000	0			No issue

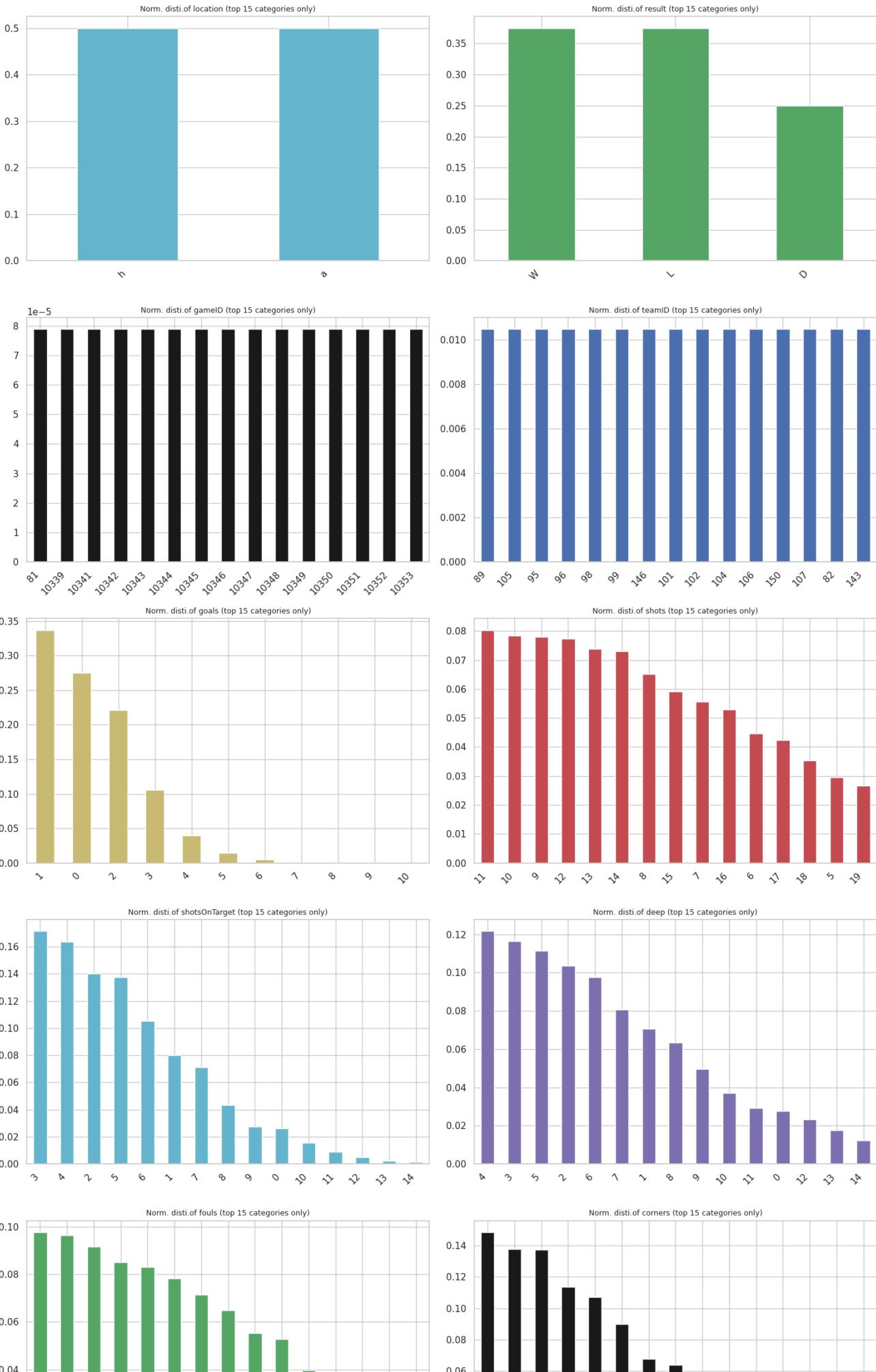
Number of All Scatter Plots = 6

Pair-wise Scatter Plot of all Continuous Variables

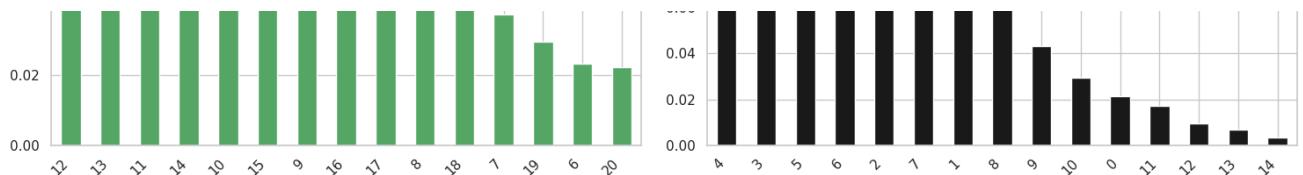




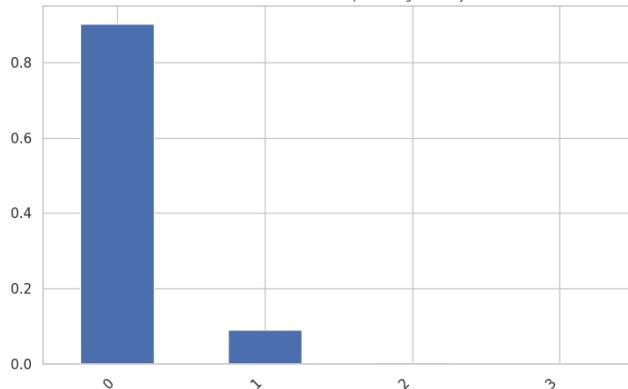
Histograms and Normalized distributions of all variables



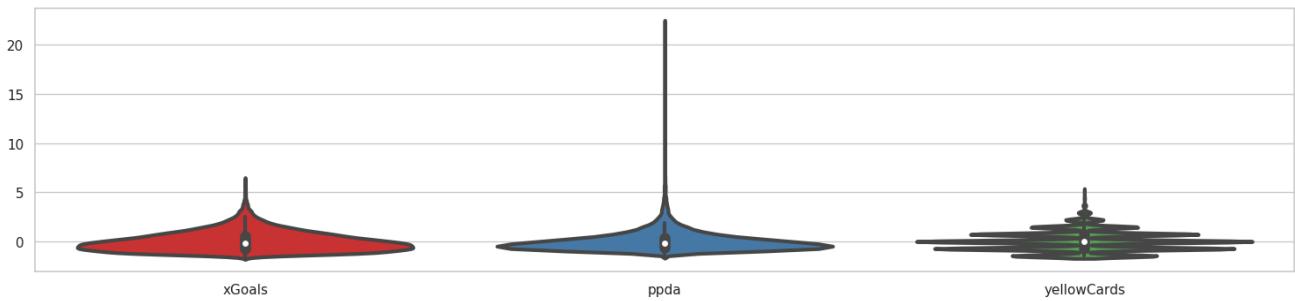
Notebook



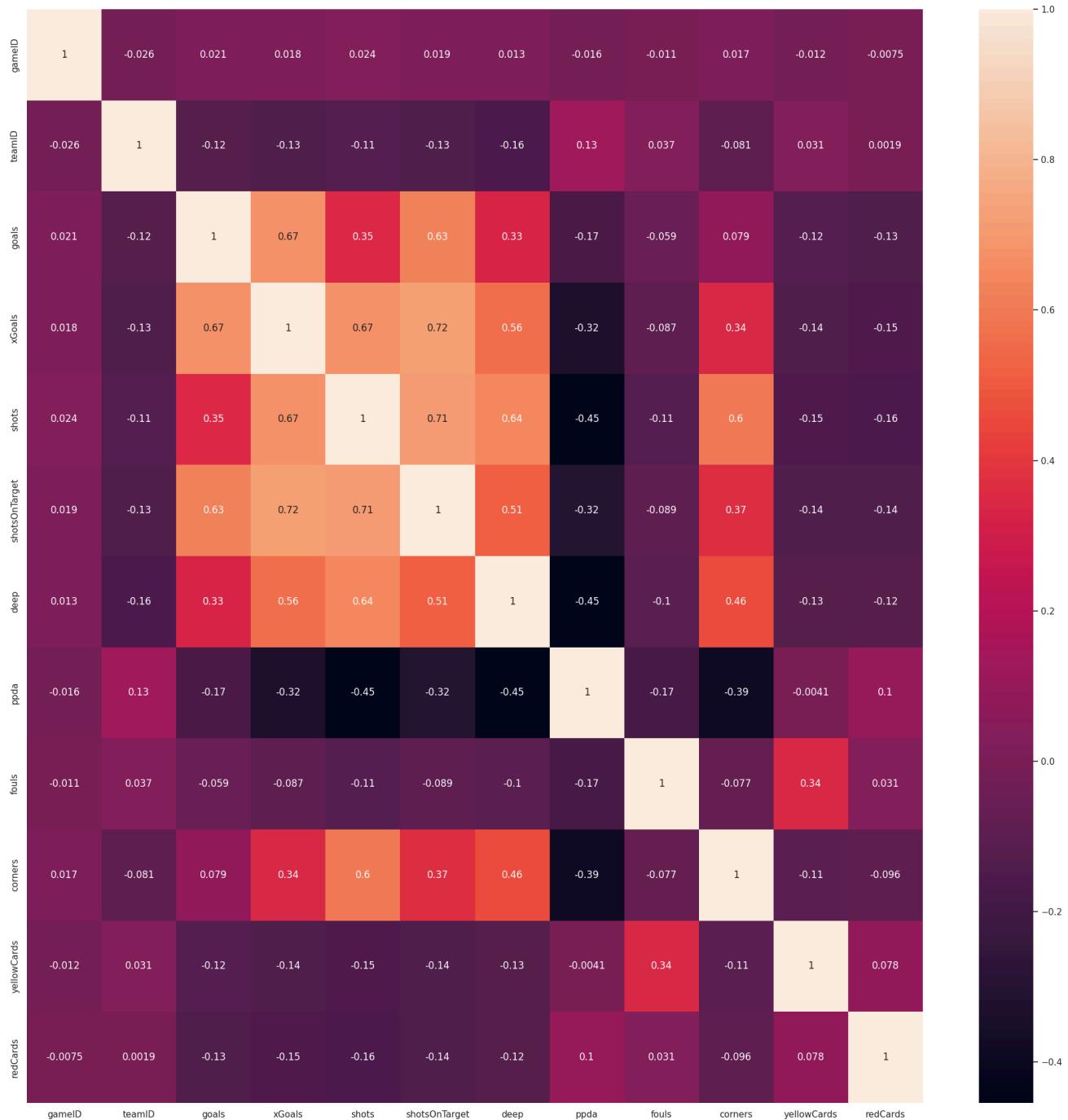
Norm. disti. of redCards (top 15 categories only)



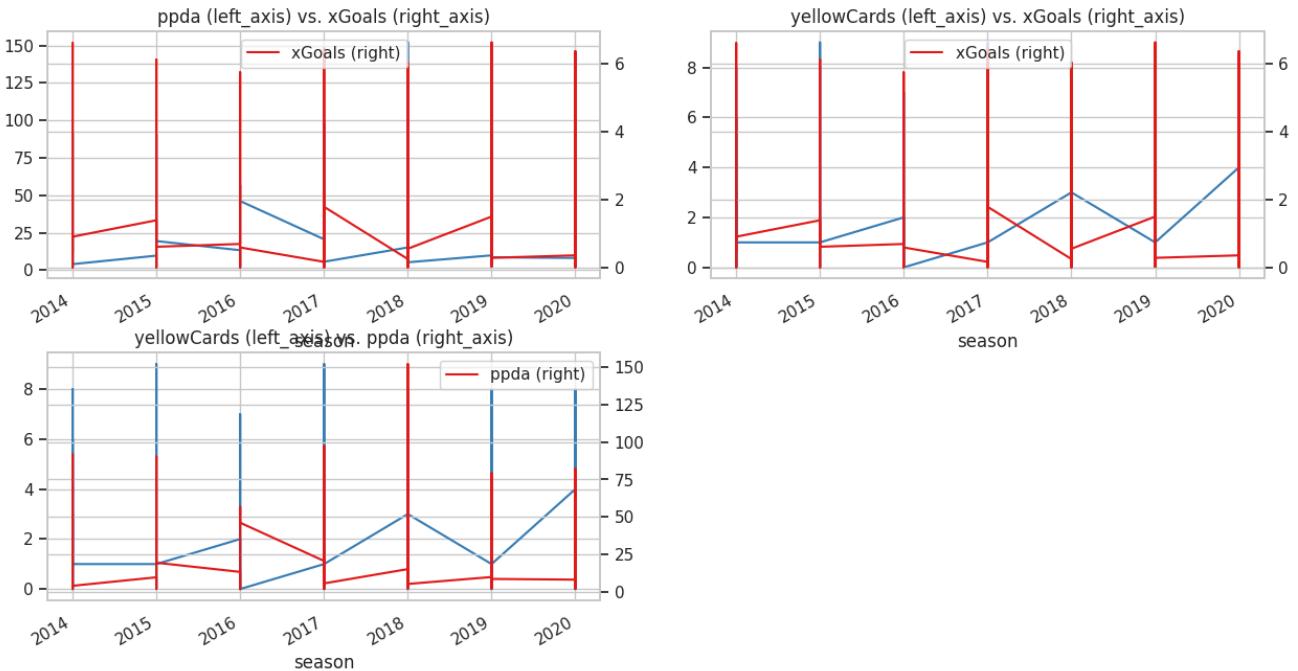
Violin Plot of all Continuous Variables



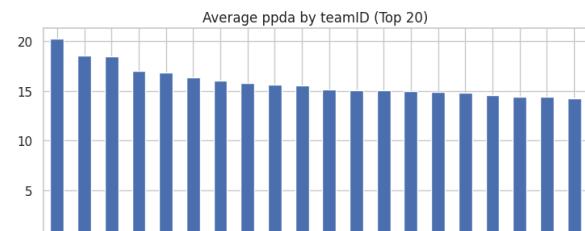
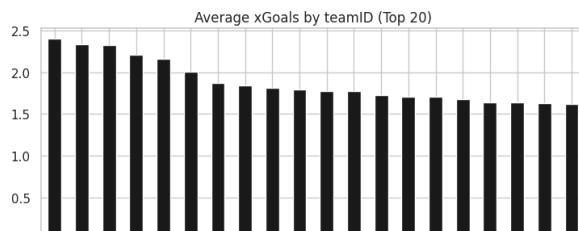
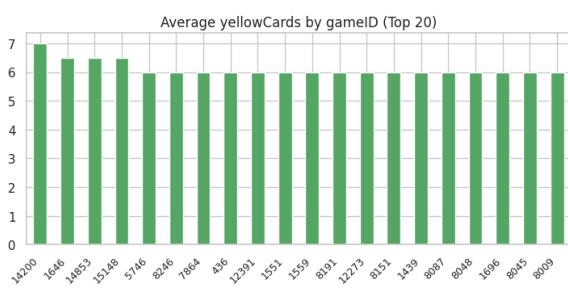
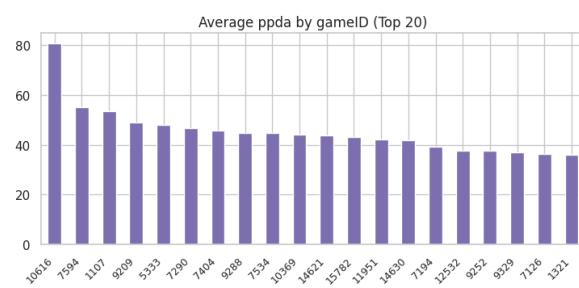
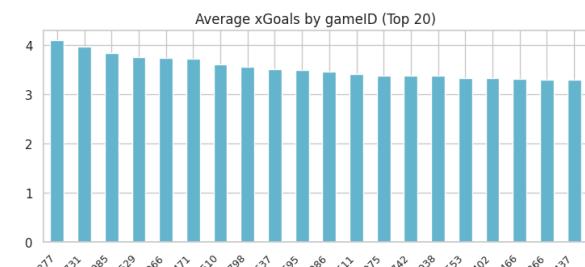
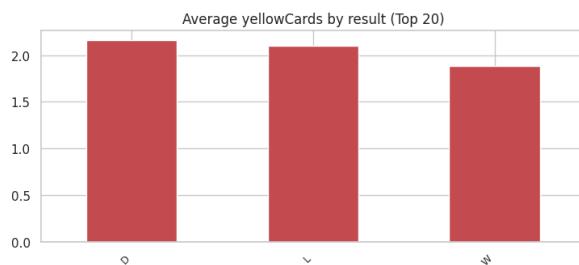
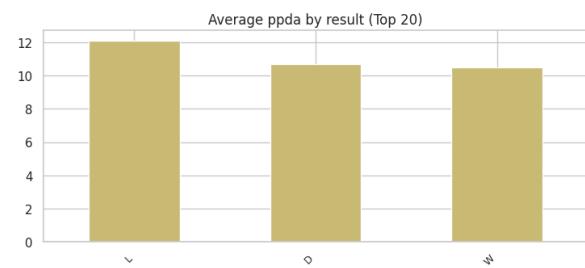
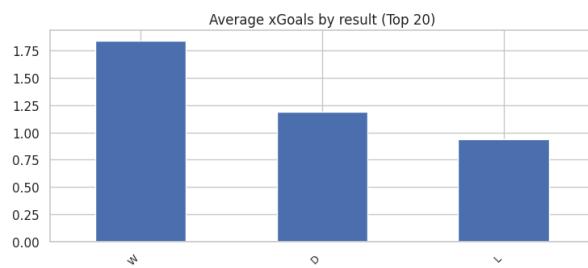
Time Series Data: Heatmap of Differenced Continuous vars including target =



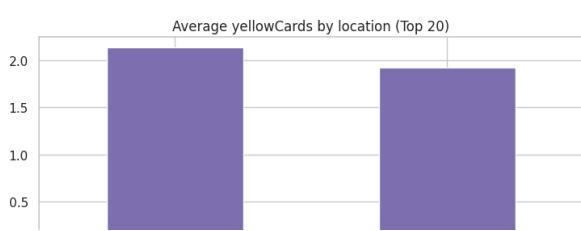
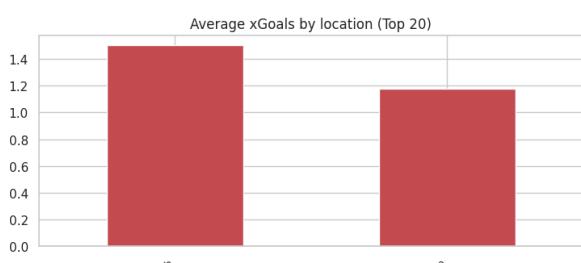
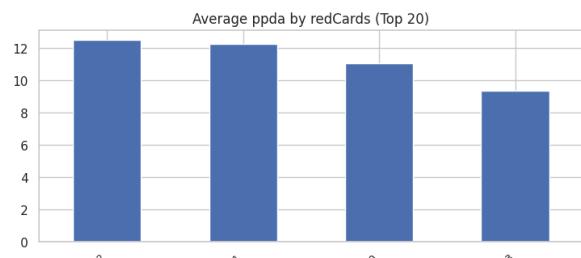
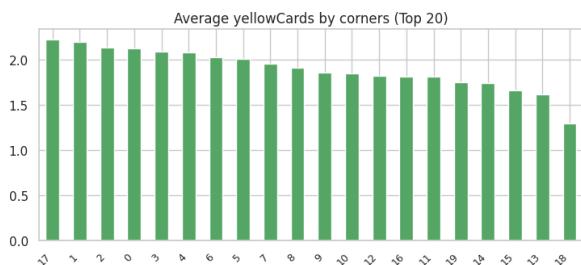
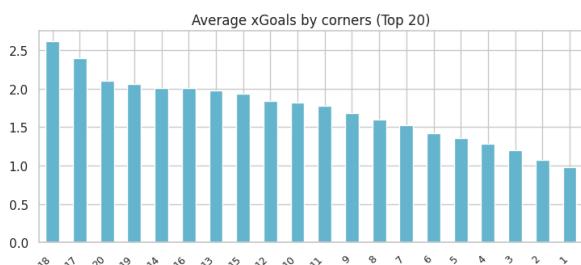
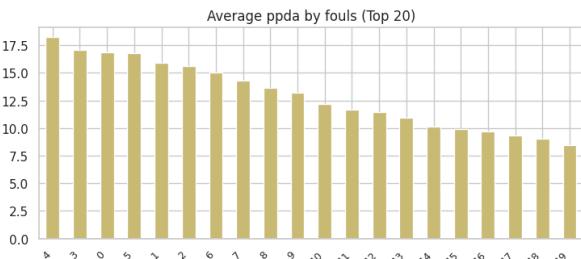
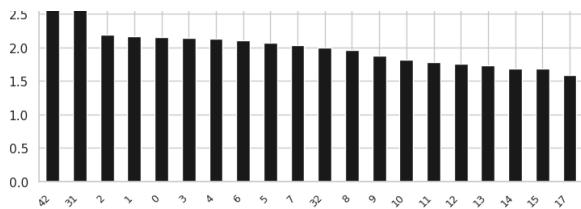
Time Series Plot by season: Pairwise Continuous Variables



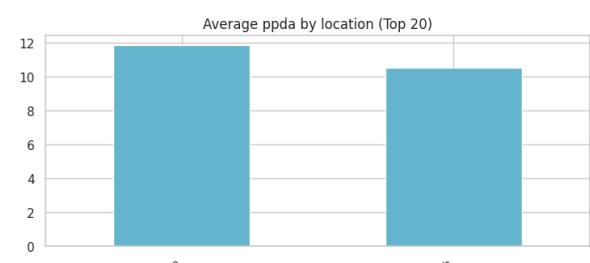
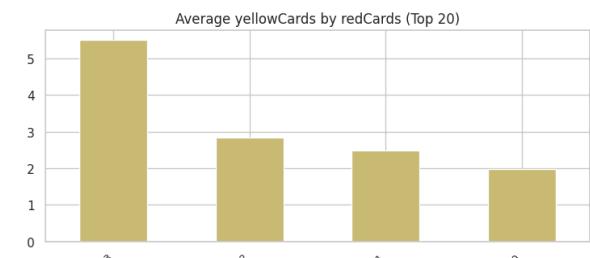
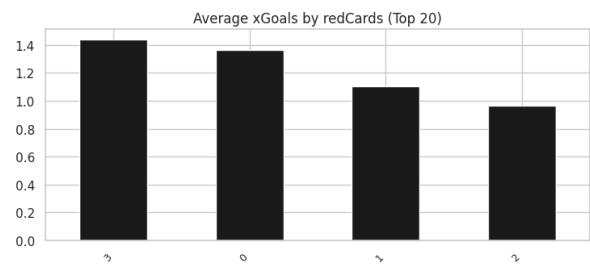
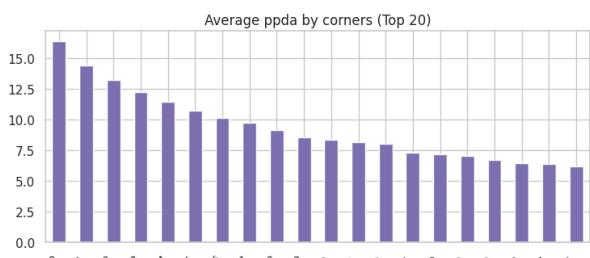
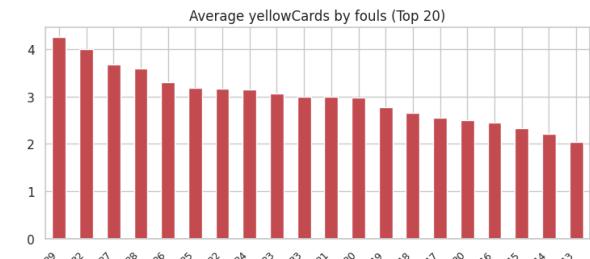
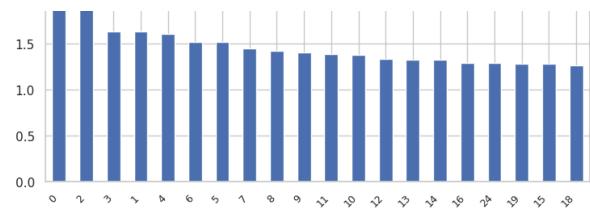
Bar plots for each Continuous by each Categorical variable



30/03/2025, 13:33



Notebook





```
[nltk_data] Downloading collection 'popular'
[nltk_data]   | Downloading package cmudict to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package cmudict is already up-to-date!
[nltk_data]   | Downloading package gazetteers to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package gazetteers is already up-to-date!
[nltk_data]   | Downloading package genesis to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package genesis is already up-to-date!
[nltk_data]   | Downloading package gutenberg to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package gutenberg is already up-to-date!
[nltk_data]   | Downloading package inaugural to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package inaugural is already up-to-date!
[nltk_data]   | Downloading package movie_reviews to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package movie_reviews is already up-to-date!
[nltk_data]   | Downloading package names to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package names is already up-to-date!
[nltk_data]   | Downloading package shakespeare to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package shakespeare is already up-to-date!
[nltk_data]   | Downloading package stopwords to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package stopwords is already up-to-date!
[nltk_data]   | Downloading package treebank to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package treebank is already up-to-date!
[nltk_data]   | Downloading package twitter_samples to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package twitter_samples is already up-to-date!
[nltk_data]   | Downloading package omw to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package omw is already up-to-date!
[nltk_data]   | Downloading package omw-1.4 to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package omw-1.4 is already up-to-date!
[nltk_data]   | Downloading package wordnet to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package wordnet is already up-to-date!
[nltk_data]   | Downloading package wordnet2021 to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package wordnet2021 is already up-to-date!
[nltk_data]   | Downloading package wordnet31 to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package wordnet31 is already up-to-date!
[nltk_data]   | Downloading package wordnet_ic to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package wordnet_ic is already up-to-date!
[nltk_data]   | Downloading package words to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package words is already up-to-date!
[nltk_data]   | Downloading package maxent_ne_chunker to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package maxent_ne_chunker is already up-to-date!
[nltk_data]   | Downloading package punkt to
```

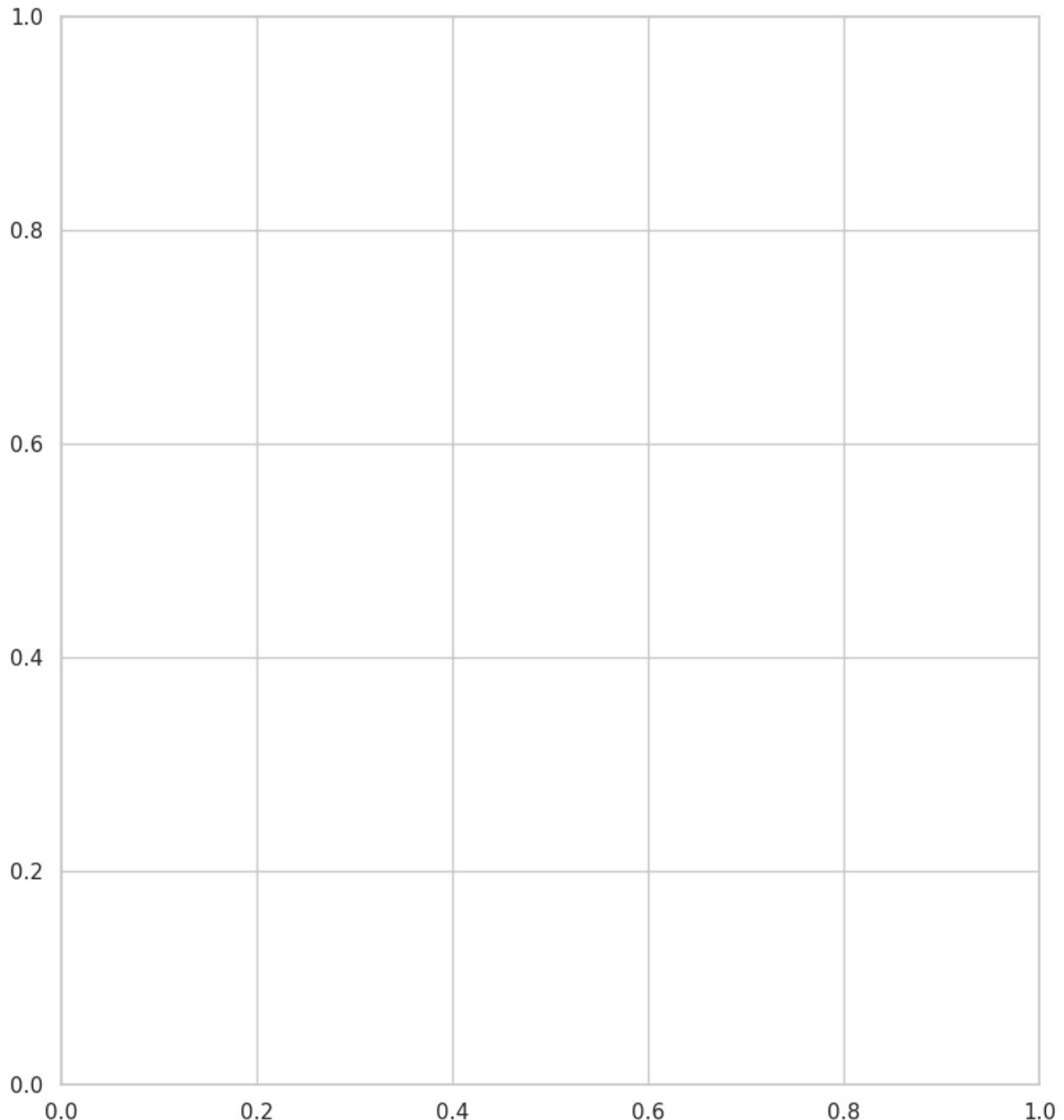
```
[nltk_data]      /home/leoadmin/nltk_data...
[nltk_data]      Package punkt is already up-to-date!
[nltk_data]      Downloading package snowball_data to
[nltk_data]          /home/leoadmin/nltk_data...
[nltk_data]      Package snowball_data is already up-to-date!
[nltk_data]      Downloading package averaged_perceptron_tagger to
[nltk_data]          /home/leoadmin/nltk_data...
[nltk_data]      Package averaged_perceptron_tagger is already up-
[nltk_data]          to-date!
[nltk_data]
[nltk_data] Done downloading collection popular
Could not draw wordcloud plot for date
All Plots done
Time to run AutoViz = 20 seconds

##### AUTO VISUALIZATION Completed #####
== Completed AutoViz on df_teamstats ==

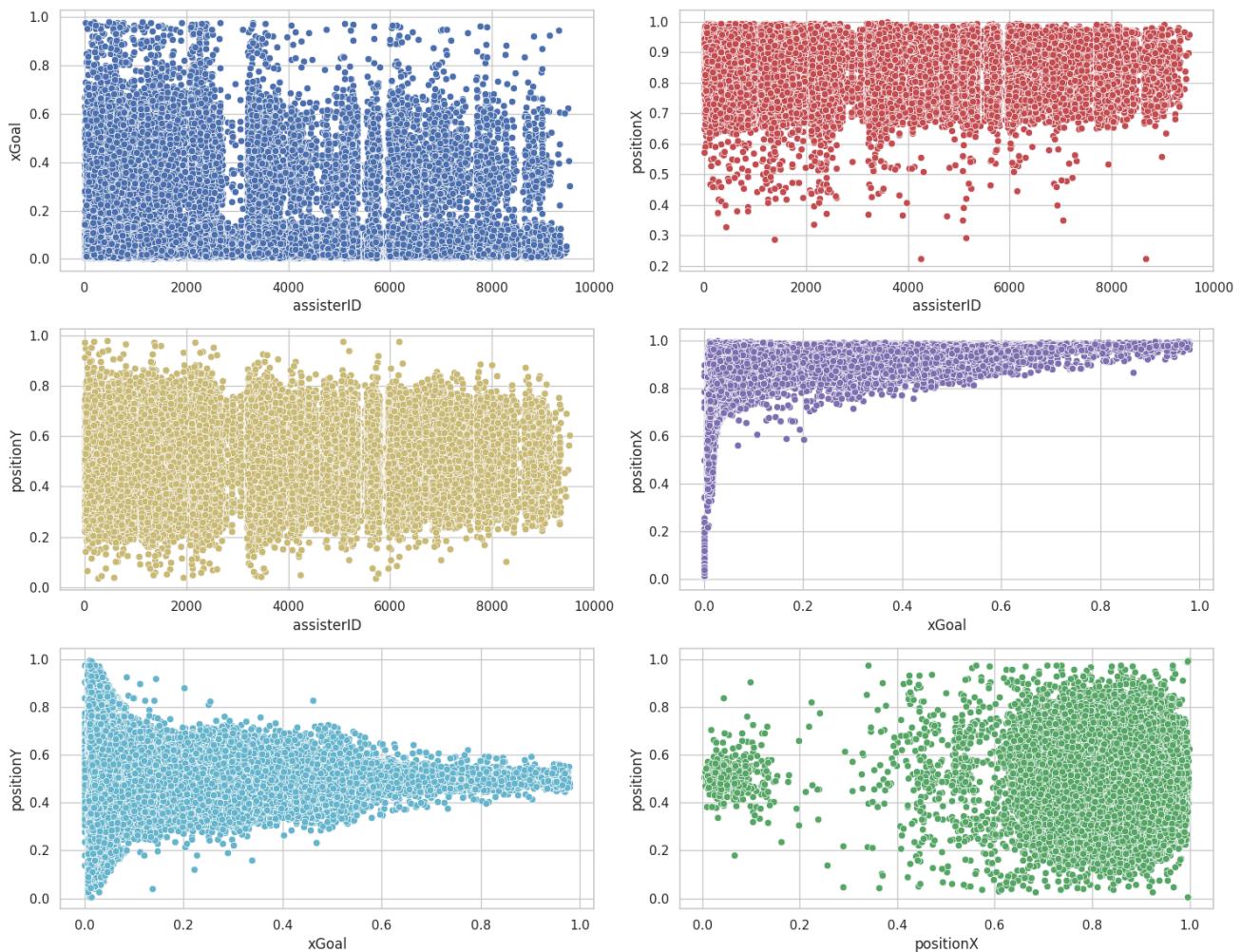
== Running AutoViz on df_shots ==
max_rows_analyzed is smaller than dataset shape 324543...
    randomly sampled 150000 rows from read CSV file
Shape of your Data Set loaded: (150000, 11)
#####
##### CLASSIFYING VARIABLES #####
#####
##### Classifying variables in data set...
Number of Numeric Columns = 4
Number of Integer-Categorical Columns = 3
Number of String-Categorical Columns = 4
Number of Factor-Categorical Columns = 0
Number of String-Boolean Columns = 0
Number of Numeric-Boolean Columns = 0
Number of Discrete String Columns = 0
Number of NLP String Columns = 0
Number of Date Time Columns = 0
Number of ID Columns = 0
Number of Columns to Delete = 0
11 Predictors classified...
    No variables removed since no ID or low-information variables found in dat
a set
Since Number of Rows in data 150000 exceeds maximum, randomly sampling 150000 rows
for EDA...
To fix these data quality issues in the dataset, import FixDQ from autoviz...
    All variables classified into correct types.
```

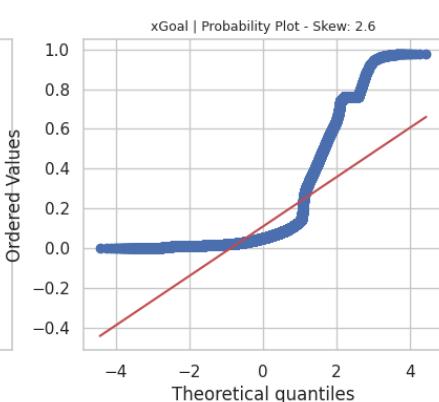
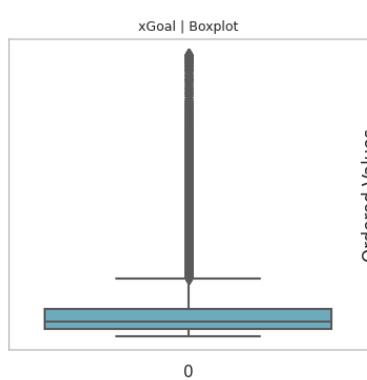
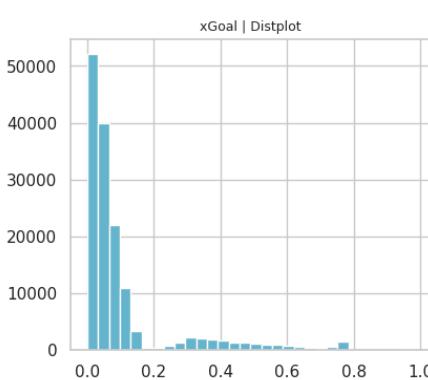
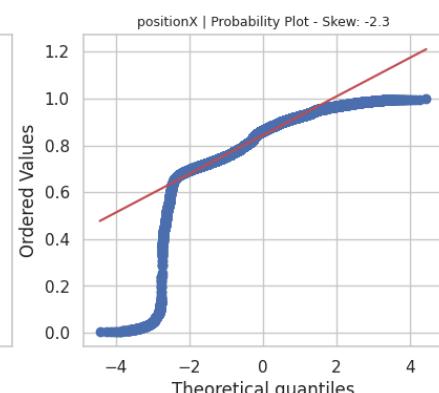
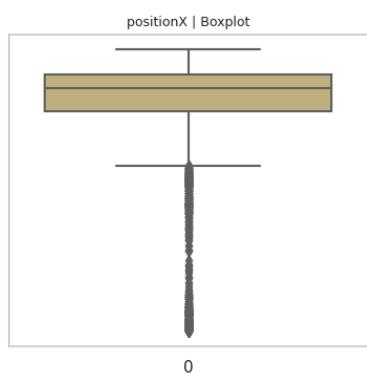
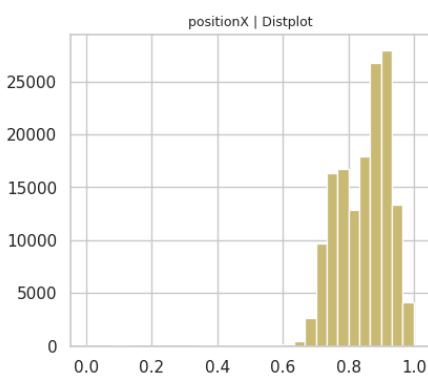
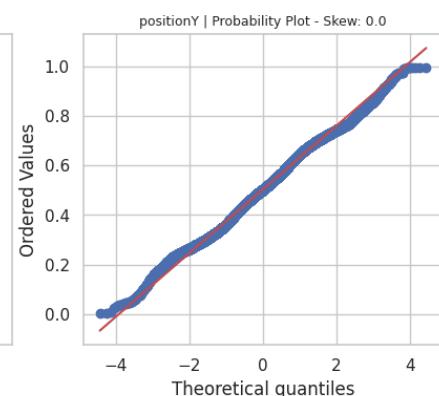
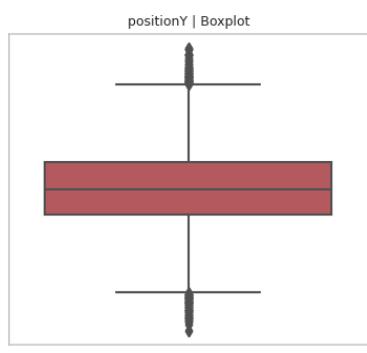
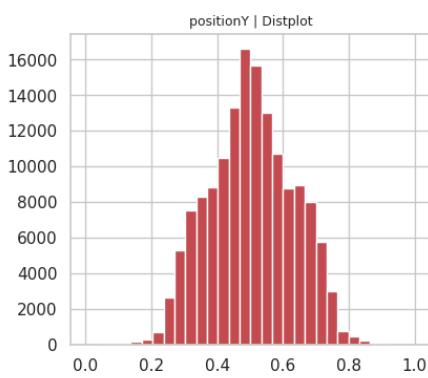
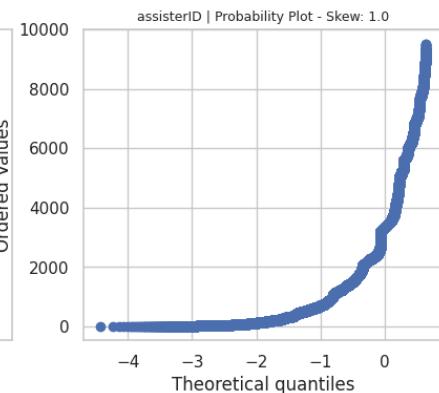
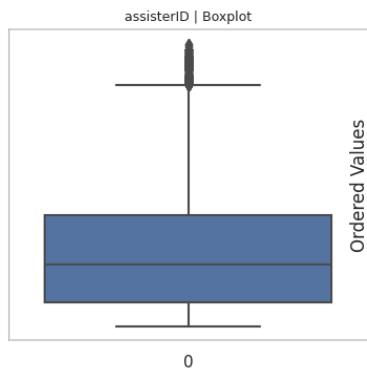
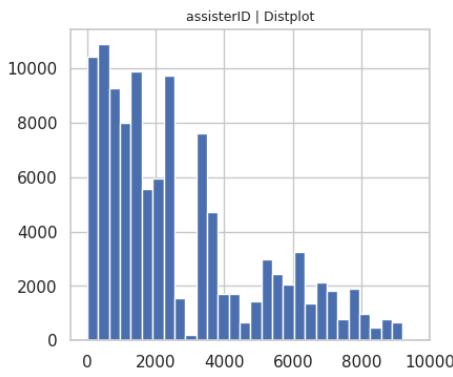
	Data Type	Missing Values%	Unique Values%	Minimum Value	Maximum Value	DQ Issue
gameID	int64	0.000000	8	81.000000	16135.000000	No issue
shooterID	int64	0.000000	3	3.000000	9566.000000	Column has 2332 outliers greater than upper bound (8225.50) or lower than lower bound(-3578.50). Cap them or remove them.
assisterID	float64	26.037333	NA	1.000000	9526.000000	39056 missing values. Impute them with mean, median, mode, or a constant value such as 123., Column has 2133 outliers greater than upper bound (8127.00) or lower than lower bound(-3553.00). Cap them or remove them.
minute	int64	0.000000	0	0.000000	103.000000	No issue
situation	object	0.000000	0			No issue
lastAction	object	0.000000	0			24 rare categories: Too many to list. Group them into a single category or drop the categories.
shotType	object	0.000000	0			1 rare categories: ['OtherBodyPart']. Group them into a single category or drop the categories.
shotResult	object	0.000000	0			1 rare categories: ['OwnGoal']. Group them into a single category or drop the categories.
xGoal	float64	0.000000	NA	0.000000	0.979344	Column has 20996 outliers greater than upper bound (0.20) or lower than lower bound(-0.08). Cap them or remove them.
positionX	float64	0.000000	NA	0.004000	0.999000	Column has 867 outliers greater than upper bound (1.10) or lower than lower bound(0.59). Cap them or remove them.
positionY	float64	0.000000	NA	0.005000	0.997000	Column has 322 outliers greater than upper bound (0.87) or lower than lower bound(0.14). Cap them or remove them.

Number of All Scatter Plots = 10

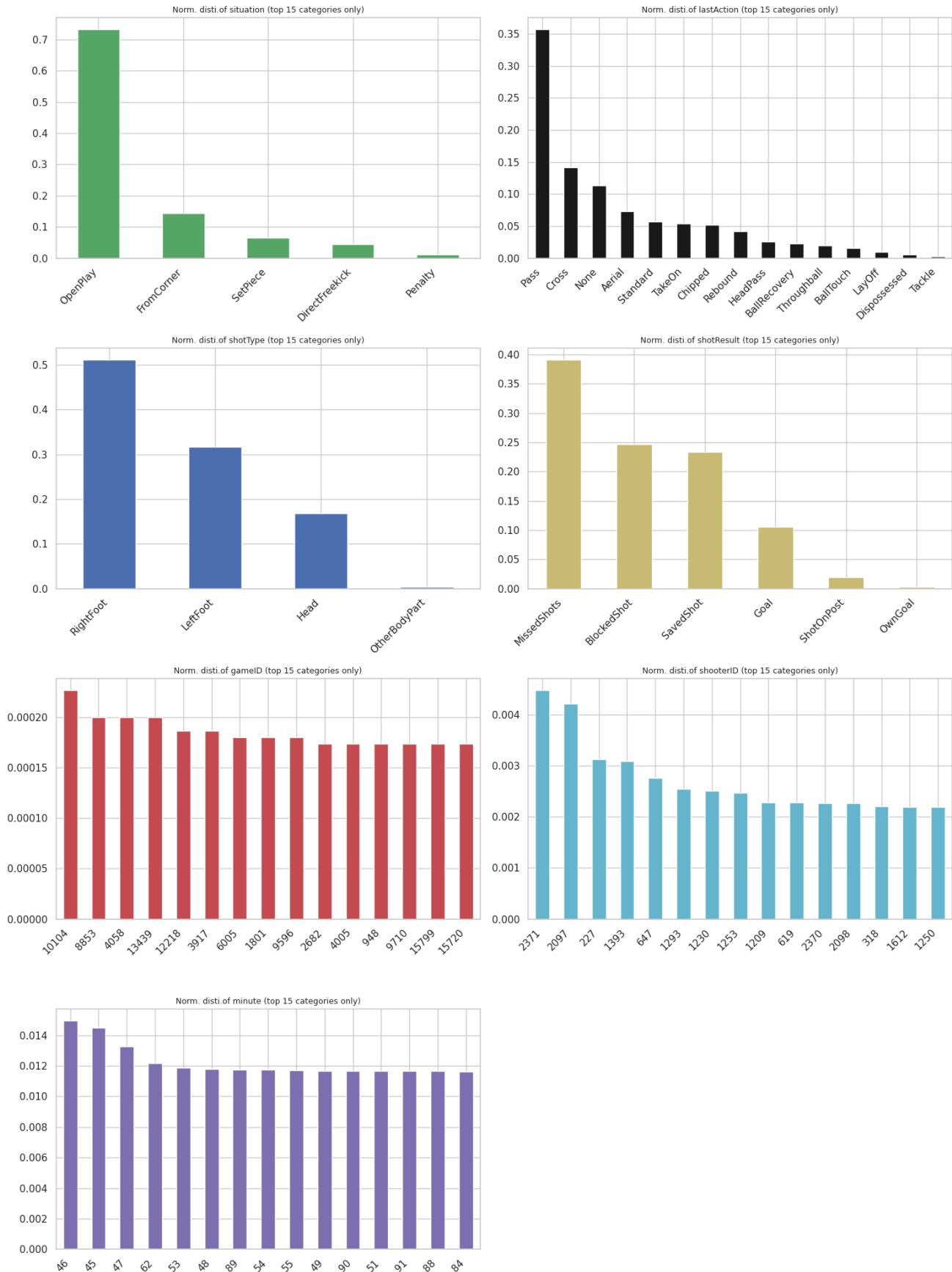


Pair-wise Scatter Plot of all Continuous Variables

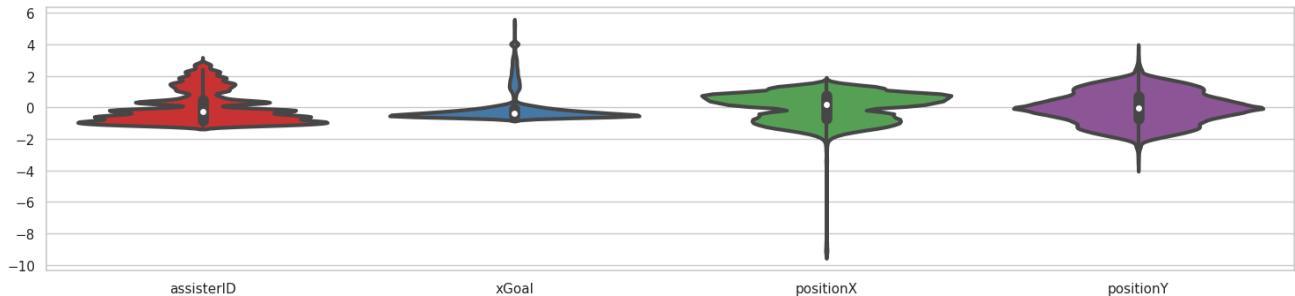




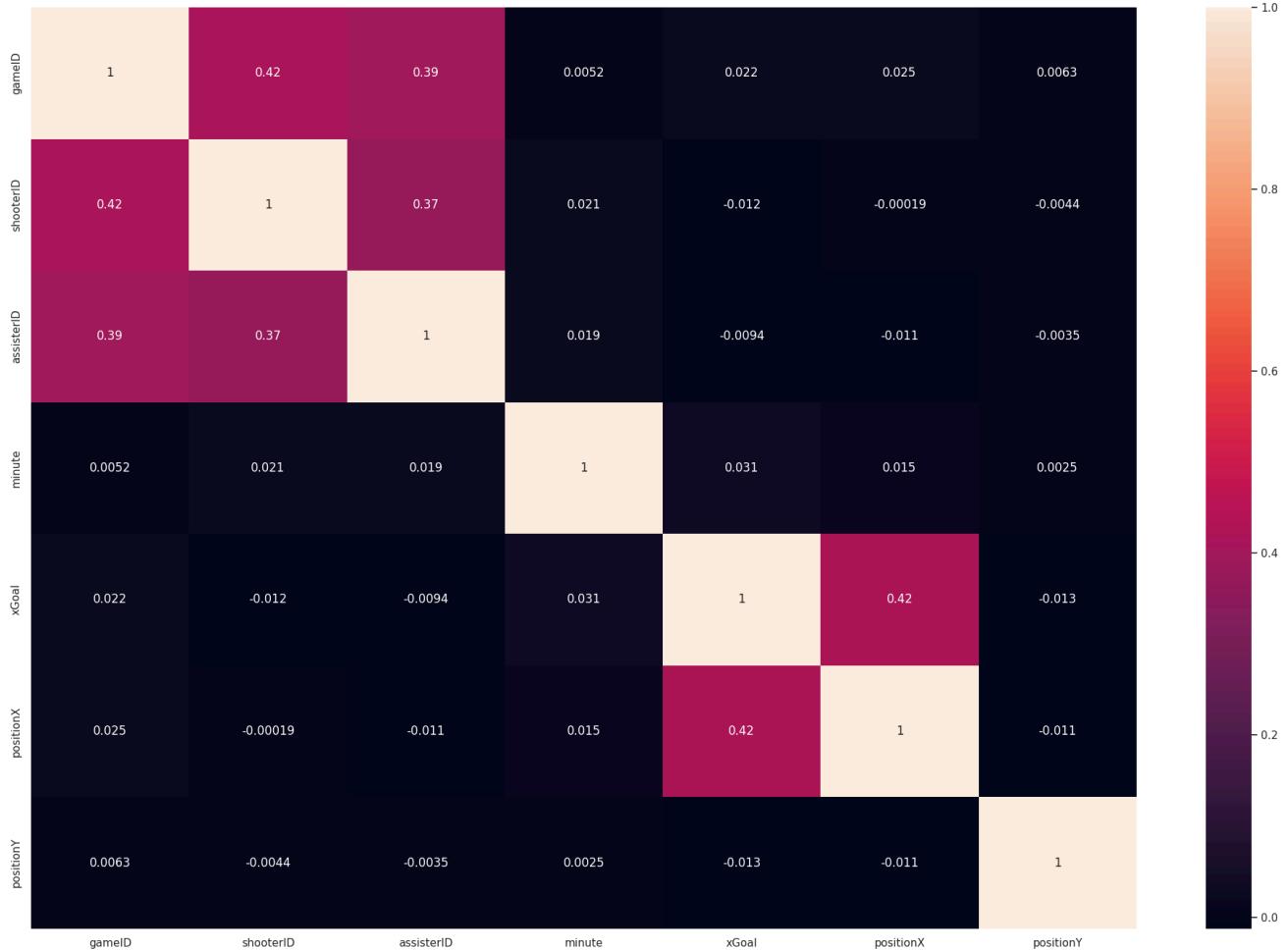
Histograms and Normalized distributions of all variables



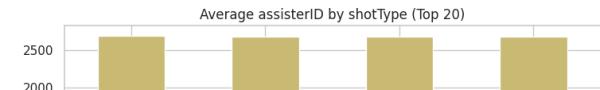
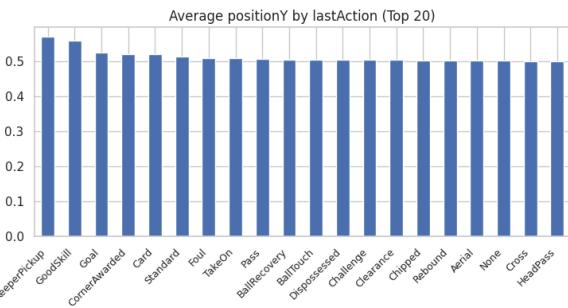
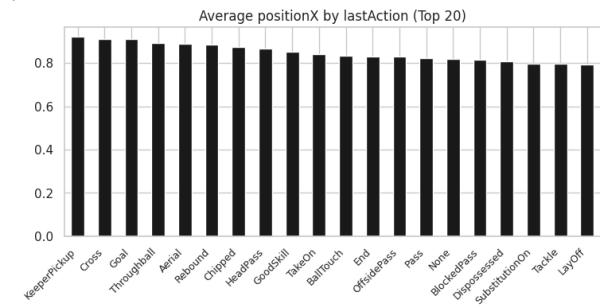
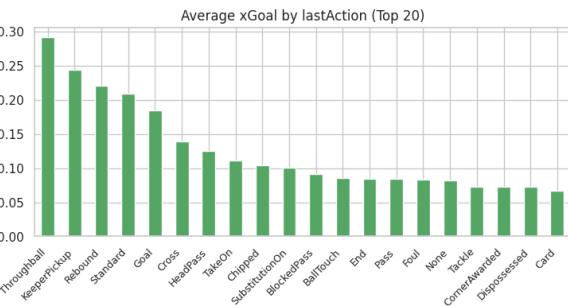
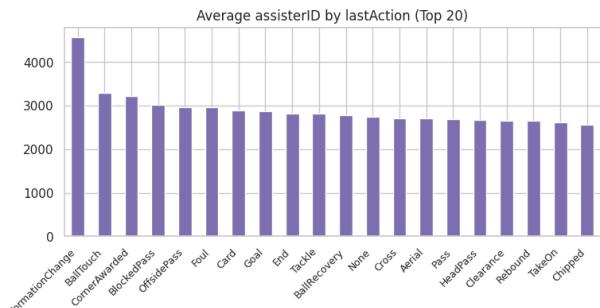
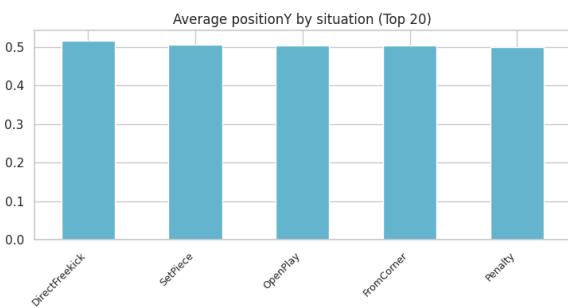
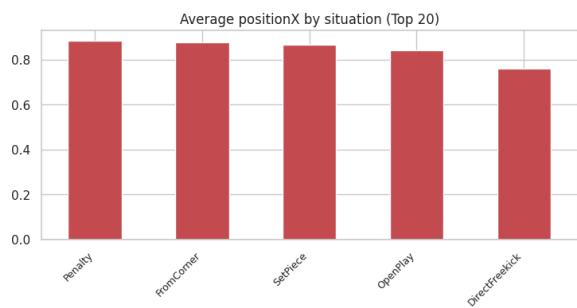
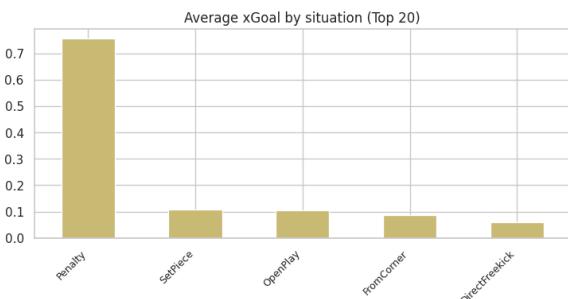
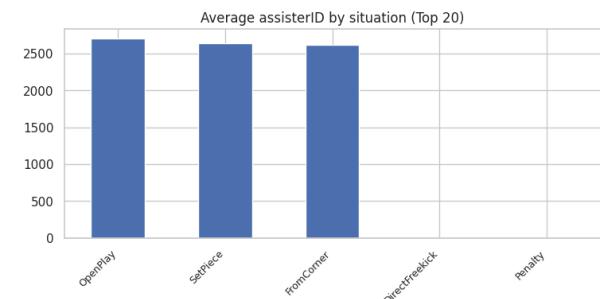
Violin Plot of all Continuous Variables

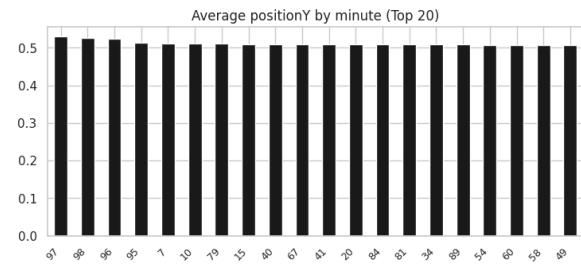
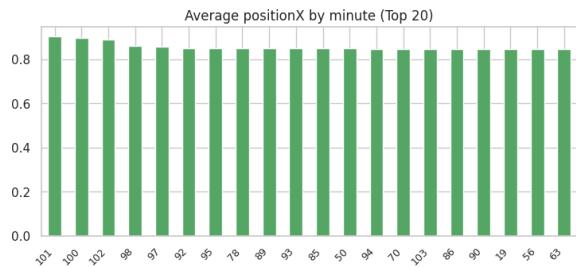
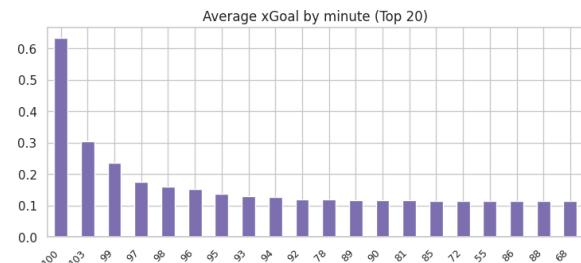
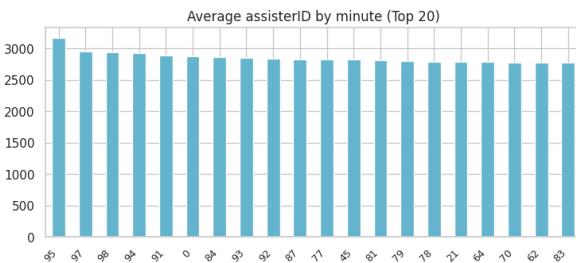
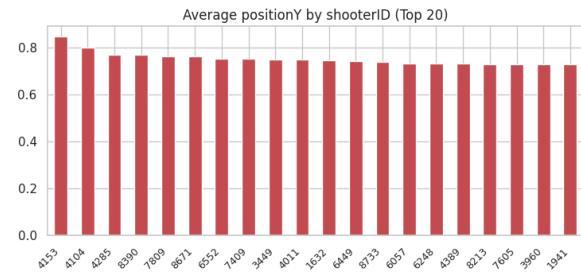
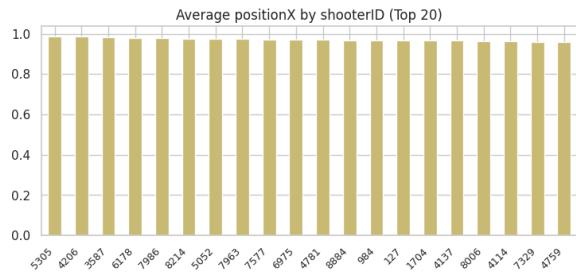


Heatmap of all Numeric Variables including target:



Bar plots for each Continuous by each Categorical variable





All Plots done

Time to run AutoViz = 16 seconds

```
##### AUTO VISUALIZATION Completed #####
== Completed AutoViz on df_shots ==
```

```
== Running AutoViz on player_shots ==
max_rows_analyzed is smaller than dataset shape 324543...
randomly sampled 150000 rows from read CSV file
```

Shape of your Data Set loaded: (150000, 11)

```
#####
##### CLASSIFYING VARIABLES #####
#####
#####
```

Classifying variables in data set...

```
Number of Numeric Columns = 4
Number of Integer-Categorical Columns = 3
Number of String-Categorical Columns = 4
Number of Factor-Categorical Columns = 0
Number of String-Boolean Columns = 0
Number of Numeric-Boolean Columns = 0
Number of Discrete String Columns = 0
Number of NLP String Columns = 0
Number of Date Time Columns = 0
Number of ID Columns = 0
Number of Columns to Delete = 0
```

11 Predictors classified...

```
No variables removed since no ID or low-information variables found in data set
```

Since Number of Rows in data 150000 exceeds maximum, randomly sampling 150000 rows for EDA...

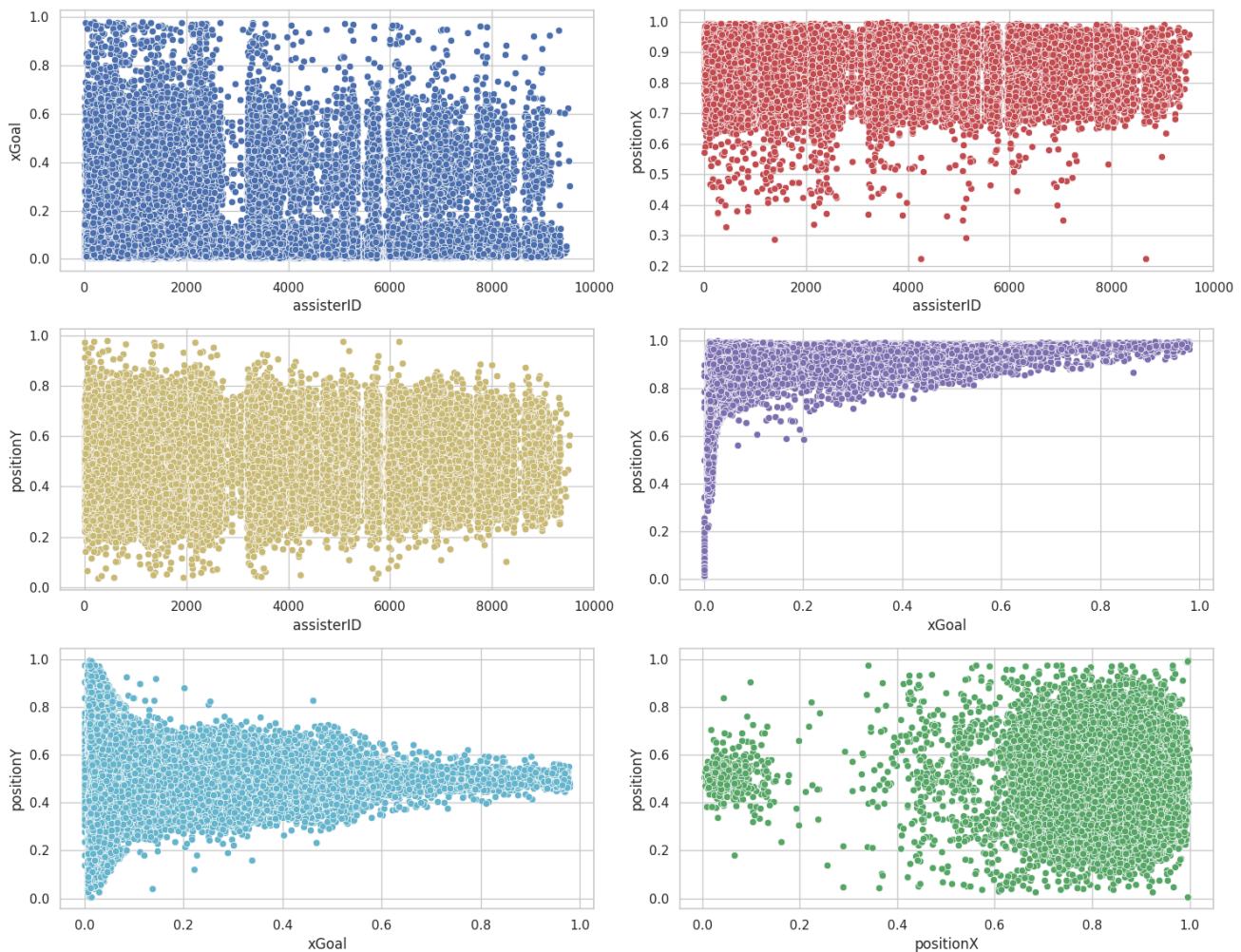
To fix these data quality issues in the dataset, import FixDQ from autoviz...

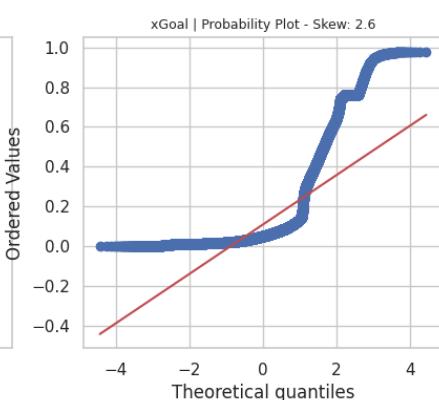
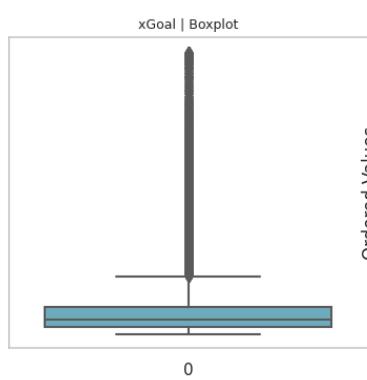
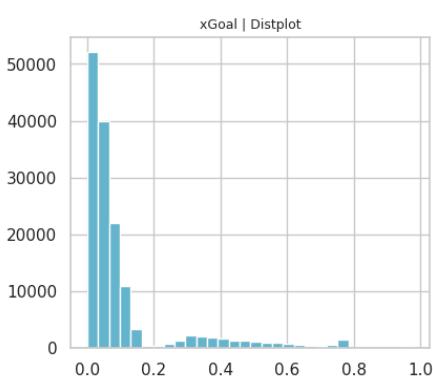
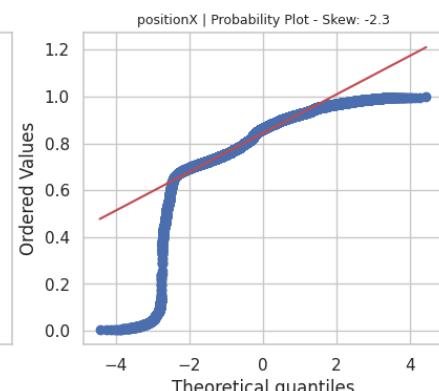
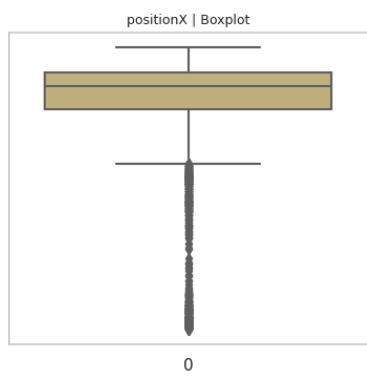
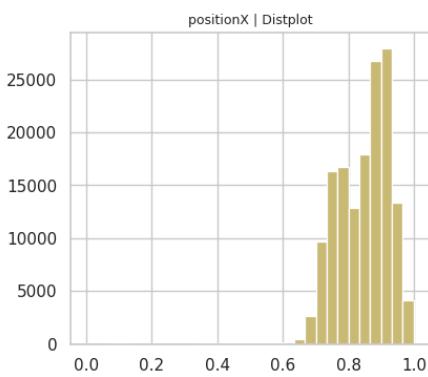
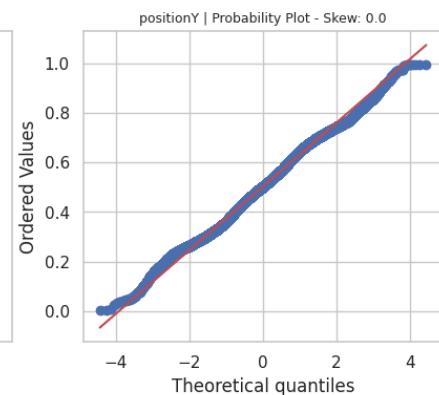
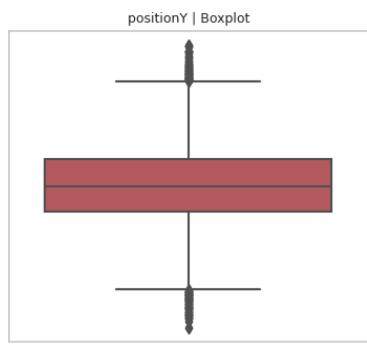
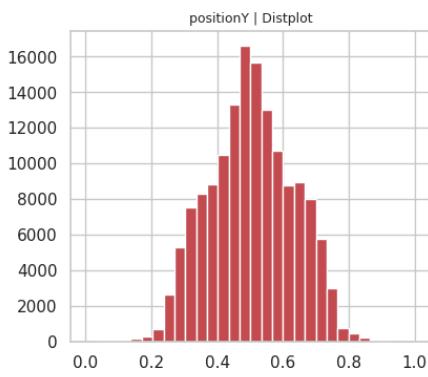
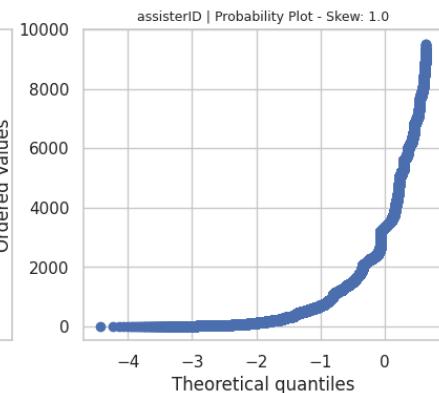
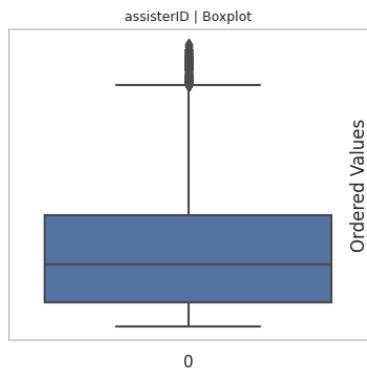
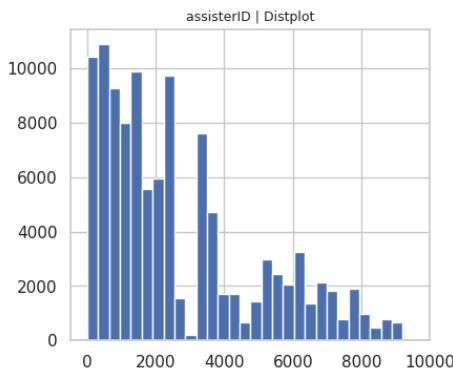
All variables classified into correct types.

	Data Type	Missing Values%	Unique Values%	Minimum Value	Maximum Value	DQ Issue
gameID	int64	0.000000	8	81.000000	16135.000000	No issue
playerID	int64	0.000000	3	3.000000	9566.000000	Column has 2332 outliers greater than upper bound (8225.50) or lower than lower bound(-3578.50). Cap them or remove them.
assisterID	float64	26.037333	NA	1.000000	9526.000000	39056 missing values. Impute them with mean, median, mode, or a constant value such as 123., Column has 2133 outliers greater than upper bound (8127.00) or lower than lower bound(-3553.00). Cap them or remove them.
minute	int64	0.000000	0	0.000000	103.000000	No issue
situation	object	0.000000	0			No issue
lastAction	object	0.000000	0			24 rare categories: Too many to list. Group them into a single category or drop the categories.
shotType	object	0.000000	0			1 rare categories: ['OtherBodyPart']. Group them into a single category or drop the categories.
shotResult	object	0.000000	0			1 rare categories: ['OwnGoal']. Group them into a single category or drop the categories.
xGoal	float64	0.000000	NA	0.000000	0.979344	Column has 20996 outliers greater than upper bound (0.20) or lower than lower bound(-0.08). Cap them or remove them.
positionX	float64	0.000000	NA	0.004000	0.999000	Column has 867 outliers greater than upper bound (1.10) or lower than lower bound(0.59). Cap them or remove them.
positionY	float64	0.000000	NA	0.005000	0.997000	Column has 322 outliers greater than upper bound (0.87) or lower than lower bound(0.14). Cap them or remove them.

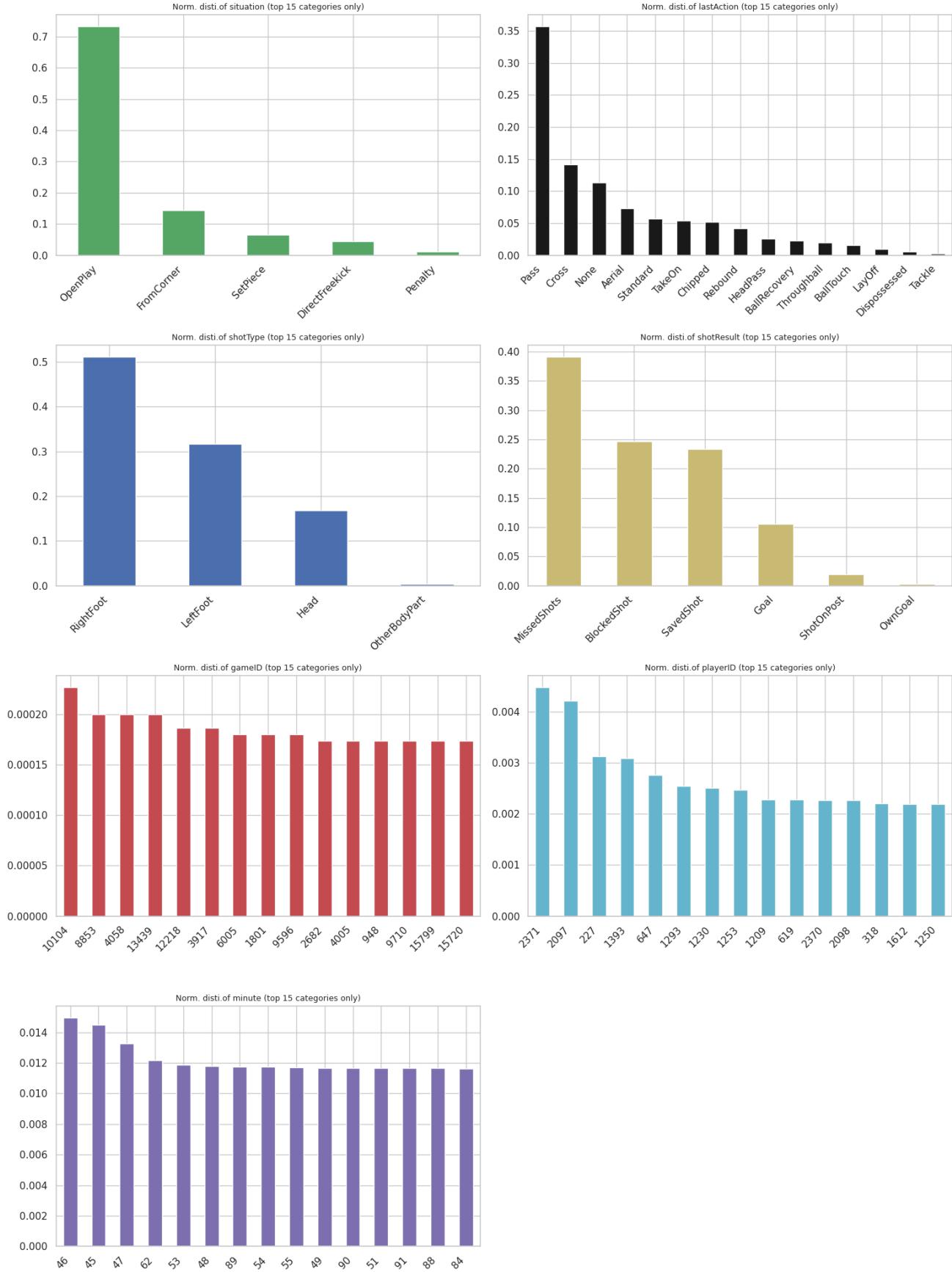
Number of All Scatter Plots = 10

Pair-wise Scatter Plot of all Continuous Variables

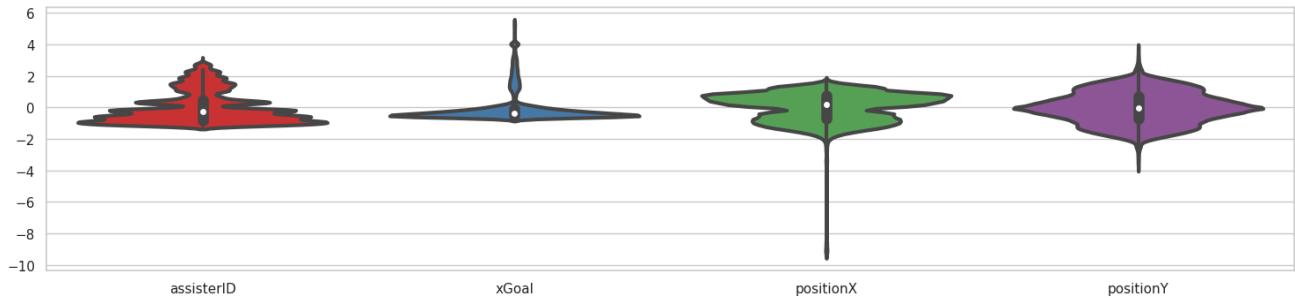




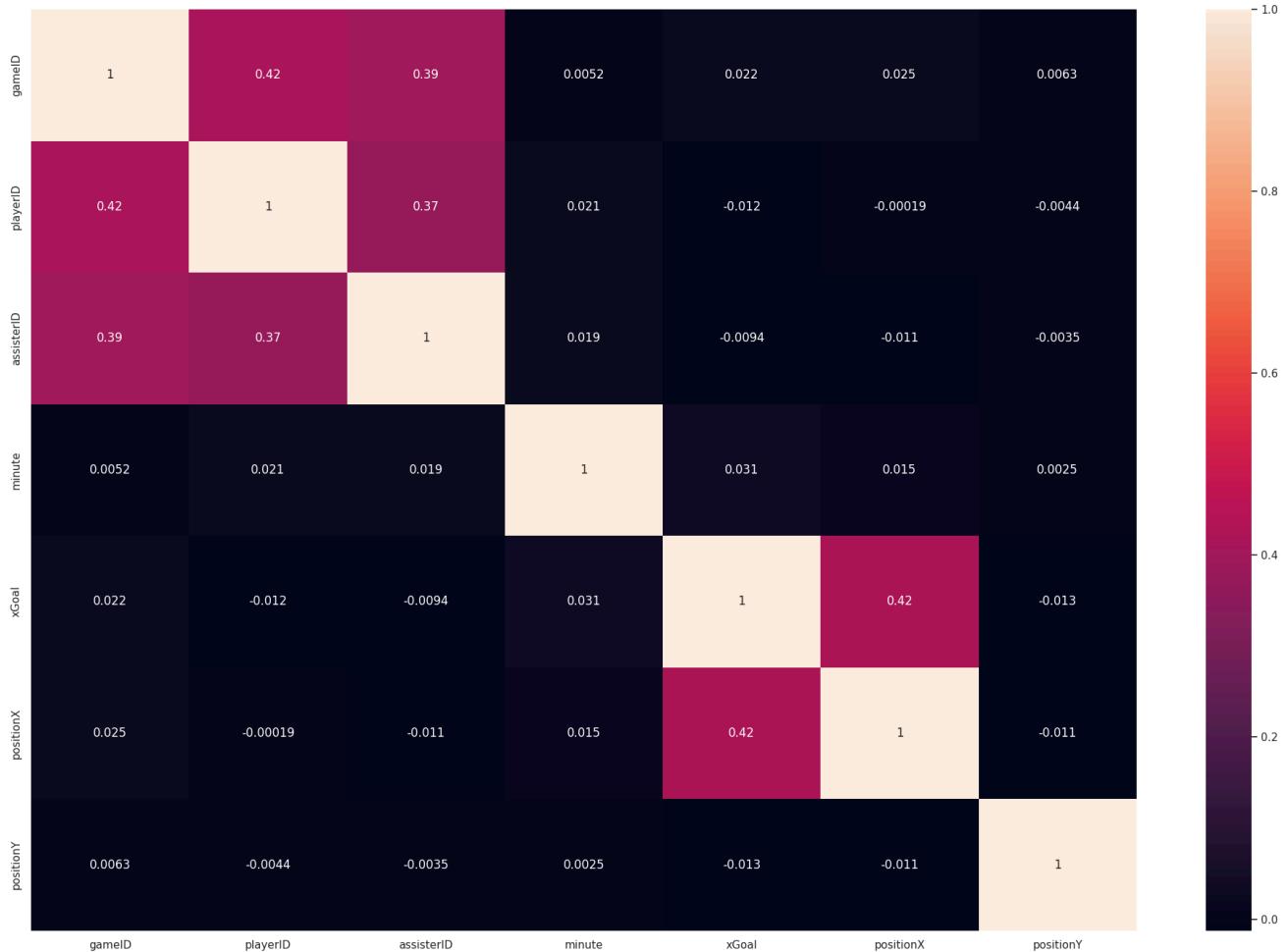
Histograms and Normalized distributions of all variables



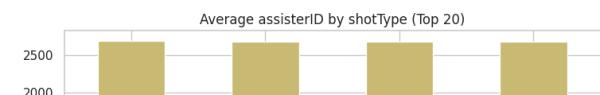
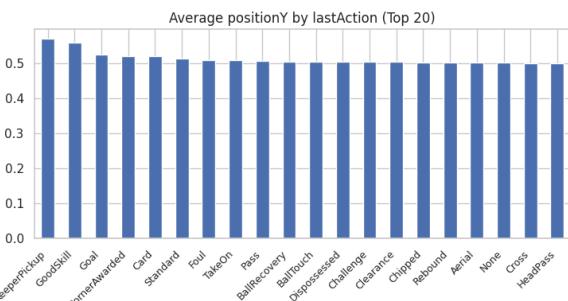
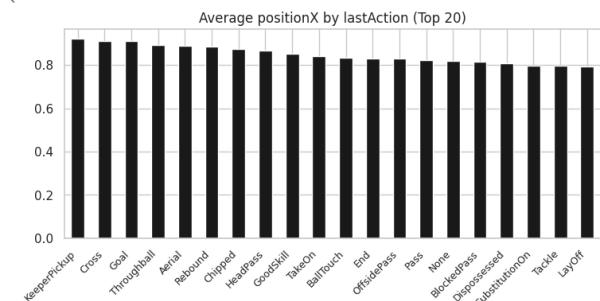
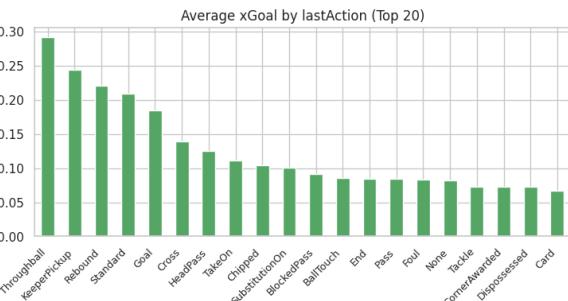
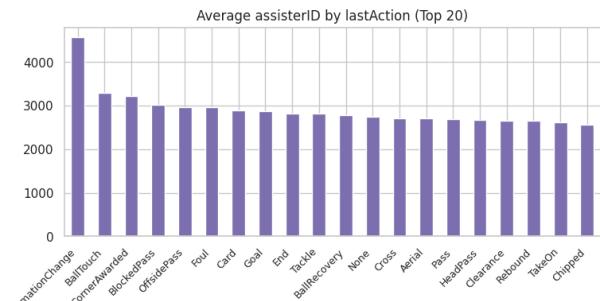
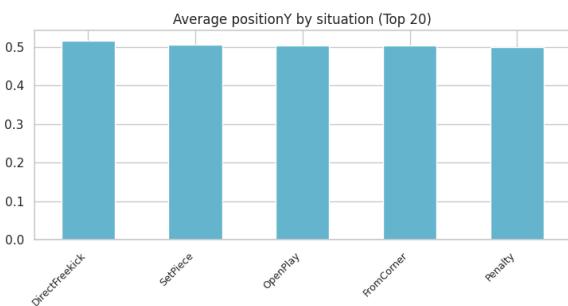
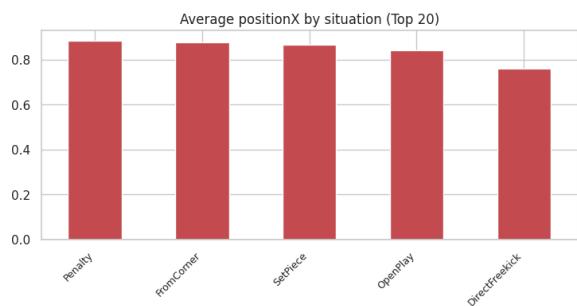
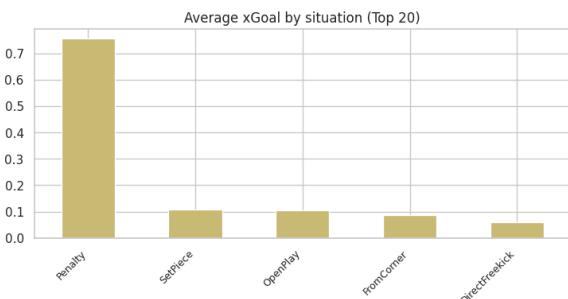
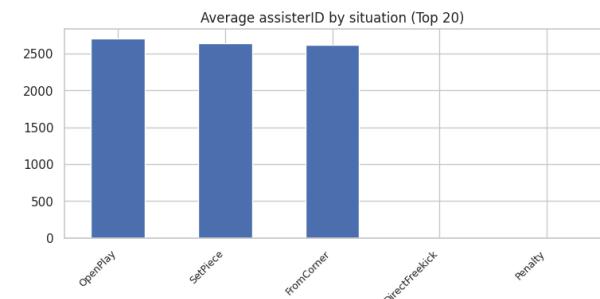
Violin Plot of all Continuous Variables

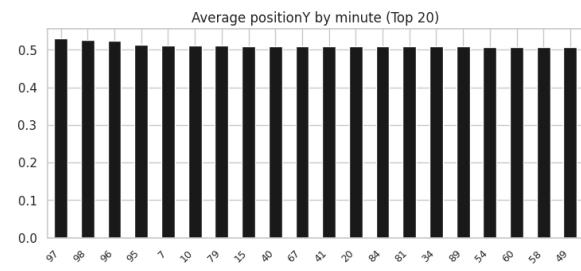
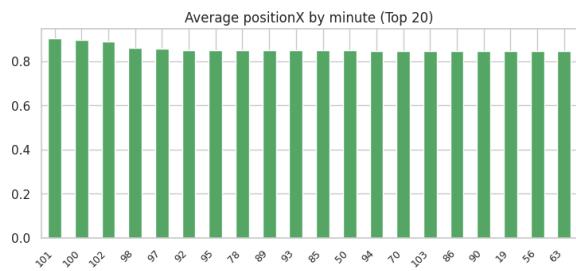
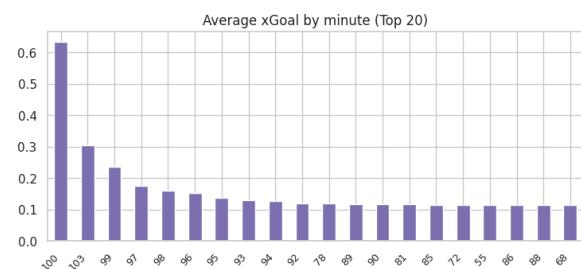
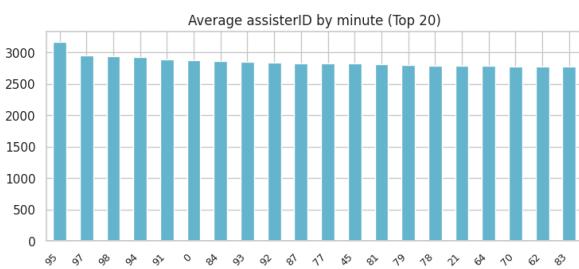
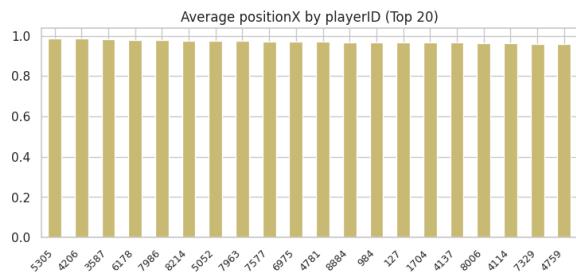


Heatmap of all Numeric Variables including target:



Bar plots for each Continuous by each Categorical variable





All Plots done

Time to run AutoViz = 16 seconds

```
##### AUTO VISUALIZATION Completed #####
== Completed AutoViz on player_shots ==
```

Skipping teamstats...

== Running AutoViz on df_combined ==

Shape of your Data Set loaded: (12680, 39)

```
#####
##### CLASSIFYING VARIABLES #####
#####
#####
```

Classifying variables in data set...

Number of Numeric Columns = 6

Number of Integer-Categorical Columns = 21

Number of String-Categorical Columns = 3

Number of Factor-Categorical Columns = 0

Number of String-Boolean Columns = 0

Number of Numeric-Boolean Columns = 0

Number of Discrete String Columns = 0

Number of NLP String Columns = 3

Number of Date Time Columns = 3

Number of ID Columns = 1

Number of Columns to Delete = 2

39 Predictors classified...

3 variable(s) removed since they were ID or low-information variables

List of variables removed: ['gameID', 'home_location', 'away_location']

6 numeric variables in data exceeds limit, taking top 30 variables

List of variables selected: ['home_xGoals', 'home_ppda', 'away_xGoals', 'away_ppda', 'home_yellowCards', 'away_yellowCards']

Total columns > 30, too numerous to print.

To fix these data quality issues in the dataset, import FixDQ from autoviz...

All variables classified into correct types.

	Data Type	Missing Values%	Unique Values%	Minimum Value	Maximum Value	DQ Issue
gameID	int64	0.000000	100	81.000000	16135.000000	Possible ID column: drop before modeling step.
leagueID	int64	0.000000	0	1.000000	5.000000	No issue
season	int64	0.000000	0	2014.000000	2020.000000	Possible date-time colum: transform before modeling step.
date	object	0.000000	53			No issue
homeTeamID	int64	0.000000	1	71.000000	262.000000	Column has 55 outliers greater than upper bound (256.00) or lower than lower bound(8.00). Cap them or remove them.
awayTeamID	int64	0.000000	1	71.000000	262.000000	Column has 55 outliers greater than upper bound (256.00) or lower than lower bound(8.00). Cap them or remove them.
homeGoals	int64	0.000000	0	0.000000	10.000000	Column has 981 outliers greater than upper bound (3.50) or lower than lower bound(-0.50). Cap them or remove them.
awayGoals	int64	0.000000	0	0.000000	9.000000	Column has 44 outliers greater than upper bound (5.00) or lower than lower bound(-3.00). Cap them or remove them.
homeGoalsHalfTime	int64	0.000000	0	0.000000	6.000000	Column has 403 outliers greater than upper bound (2.50) or lower than lower bound(-1.50). Cap them or remove them.
awayGoalsHalfTime	int64	0.000000	0	0.000000	5.000000	Column has 225 outliers greater than upper bound (2.50) or lower than lower bound(-1.50). Cap them or remove them.

	Data Type	Missing Values%	Unique Values%	Minimum Value	Maximum Value	DQ Issue
home_season	int64	0.000000	0	2014.000000	2020.000000	Possible date-time colum: transform before modeling step.
home_date	object	0.000000	53			No issue
home_location	object	0.000000	0			Possible Zero-variance or low information colum: drop before modeling step.
home_goals	int64	0.000000	0	0.000000	10.000000	Column has 981 outliers greater than upper bound (3.50) or lower than lower bound(-0.50). Cap them or remove them., Column has a high correlation with ['homeGoals']. Consider dropping one of them.
home_xGoals	float64	0.000000	NA	0.000000	6.630490	Column has 276 outliers greater than upper bound (3.76) or lower than lower bound(-0.92). Cap them or remove them.
home_shots	int64	0.000000	0	0.000000	47.000000	Column has 167 outliers greater than upper bound (27.50) or lower than lower bound(-0.50). Cap them or remove them.
home_shotsOnTarget	int64	0.000000	0	0.000000	18.000000	Column has 351 outliers greater than upper bound (10.50) or lower than lower bound(-1.50). Cap them or remove them.
home_deep	int64	0.000000	0	0.000000	42.000000	Column has 227 outliers greater than upper bound (18.00) or lower than lower bound(-6.00). Cap them or remove them.
home_ppda	float64	0.000000	NA	1.897400	97.333300	Column has 540 outliers greater than upper bound (21.31) or lower than lower bound(-1.88). Cap

	Data Type	Missing Values%	Unique Values%	Minimum Value	Maximum Value	DQ Issue
						them or remove them.
home_fouls	int64	0.000000	0	0.000000	33.000000	Column has 241 outliers greater than upper bound (22.50) or lower than lower bound(2.50). Cap them or remove them.
home_corners	int64	0.000000	0	0.000000	20.000000	Column has 143 outliers greater than upper bound (13.00) or lower than lower bound(-3.00). Cap them or remove them.
home_yellowCards	float64	0.007886	NA	0.000000	8.000000	1 missing values. Impute them with mean, median, mode, or a constant value such as 123., Column has 37 outliers greater than upper bound (6.00) or lower than lower bound(-2.00). Cap them or remove them.
home_redCards	int64	0.000000	0	0.000000	3.000000	Column has 1078 outliers greater than upper bound (0.00) or lower than lower bound(0.00). Cap them or remove them.
home_result	object	0.000000	0			No issue
away_season	int64	0.000000	0	2014.000000	2020.000000	Possible date-time colum: transform before modeling step.
away_date	object	0.000000	53			No issue
away_location	object	0.000000	0			Possible Zero-variance or low information colum: drop before modeling step.
away_goals	int64	0.000000	0	0.000000	9.000000	Column has 44 outliers greater than upper bound (5.00) or lower than lower bound(-3.00). Cap them or remove them., Column has a high correlation with ['awayGoals'].

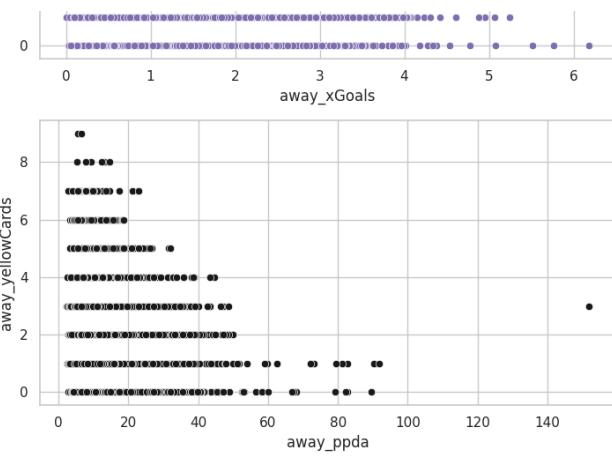
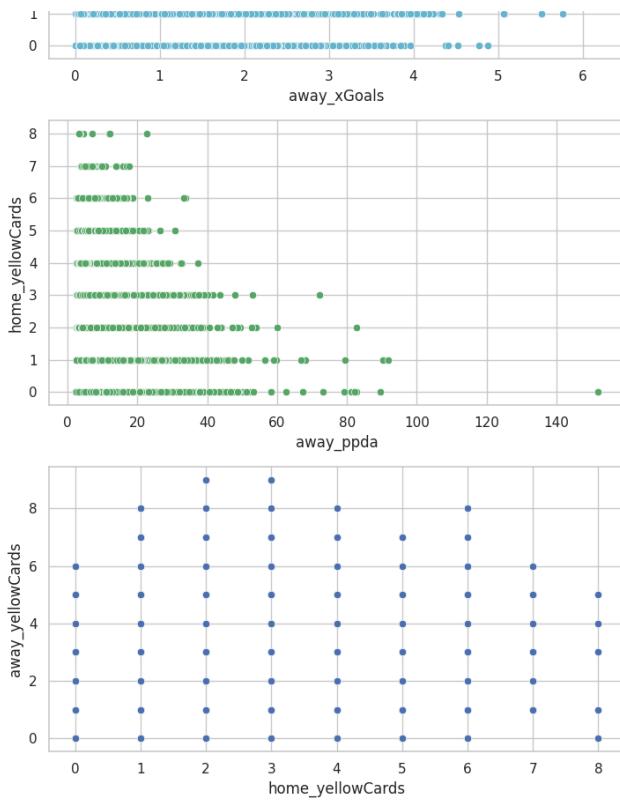
	Data Type	Missing Values%	Unique Values%	Minimum Value	Maximum Value	DQ Issue
						Consider dropping one of them.
away_xGoals	float64	0.000000	NA	0.000000	6.186960	Column has 312 outliers greater than upper bound (3.10) or lower than lower bound(-0.91). Cap them or remove them.
away_shots	int64	0.000000	0	0.000000	39.000000	Column has 161 outliers greater than upper bound (23.00) or lower than lower bound(-1.00). Cap them or remove them.
away_shotsOnTarget	int64	0.000000	0	0.000000	15.000000	Column has 233 outliers greater than upper bound (9.50) or lower than lower bound(-2.50). Cap them or remove them.
away_deep	int64	0.000000	0	0.000000	28.000000	Column has 423 outliers greater than upper bound (13.00) or lower than lower bound(-3.00). Cap them or remove them.
away_ppda	float64	0.000000	NA	2.122000	152.000000	Column has 606 outliers greater than upper bound (24.25) or lower than lower bound(-2.58). Cap them or remove them.
away_fouls	int64	0.000000	0	0.000000	32.000000	Column has 81 outliers greater than upper bound (25.00) or lower than lower bound(1.00). Cap them or remove them.
away_corners	int64	0.000000	0	0.000000	19.000000	Column has 283 outliers greater than upper bound (10.50) or lower than lower bound(-1.50). Cap them or remove them.
away_yellowCards	float64	0.000000	NA	0.000000	9.000000	Column has 43 outliers greater than upper bound (6.00) or lower than lower

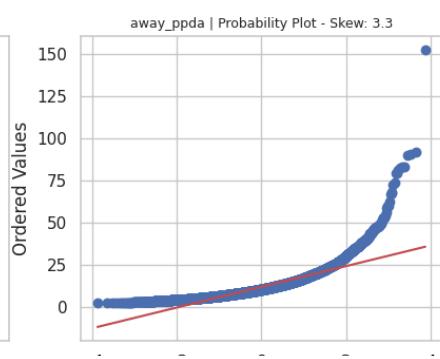
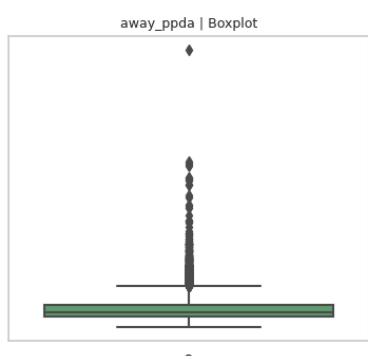
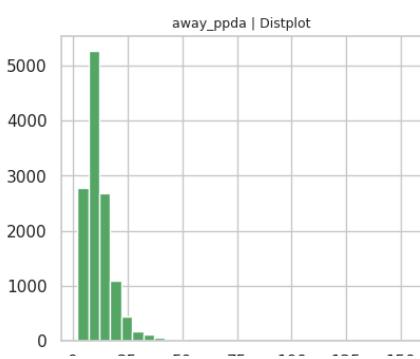
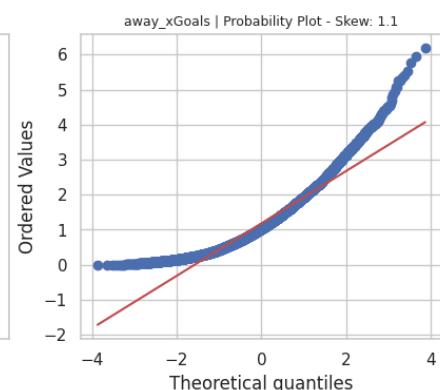
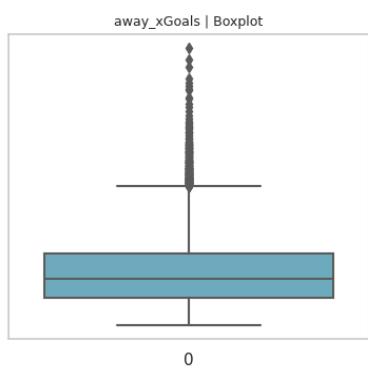
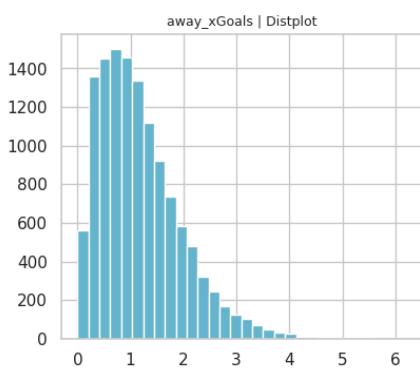
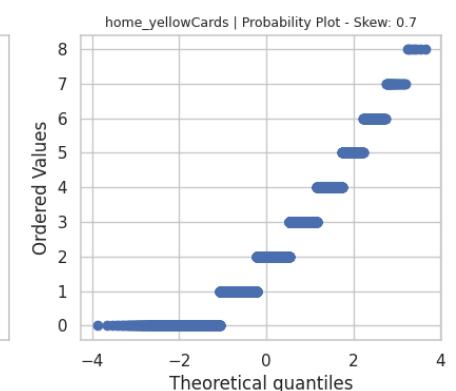
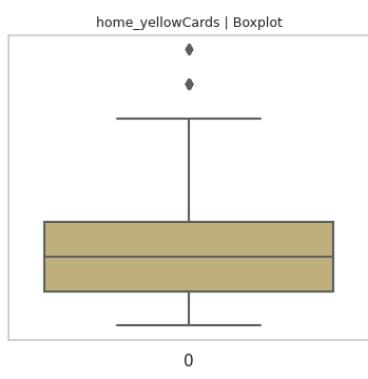
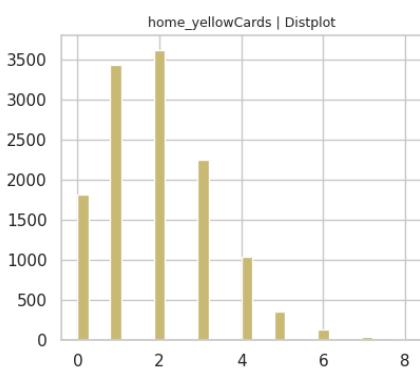
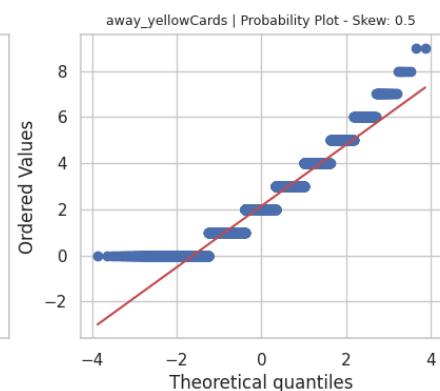
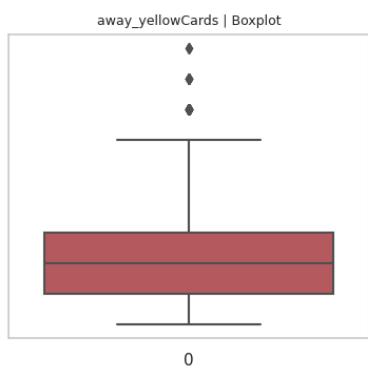
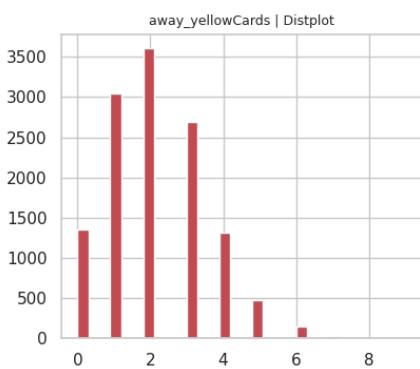
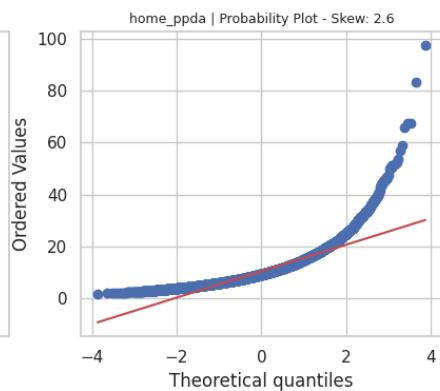
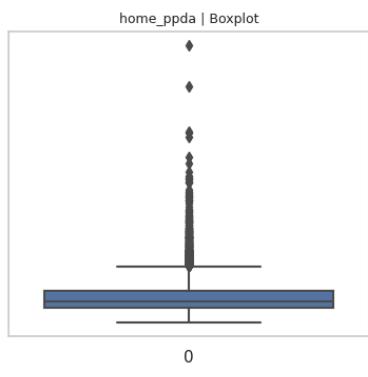
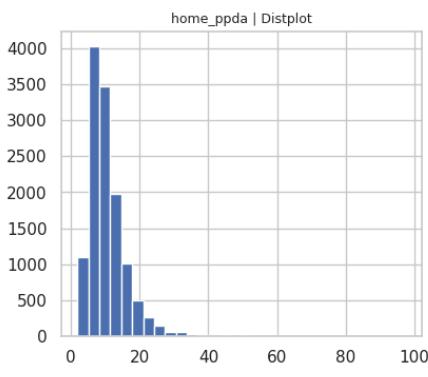
	Data Type	Missing Values%	Unique Values%	Minimum Value	Maximum Value	DQ Issue
						bound(-2.00). Cap them or remove them.
away_redCards	int64	0.000000	0	0.000000	3.000000	Column has 1396 outliers greater than upper bound (0.00) or lower than lower bound(0.00). Cap them or remove them.
away_result	object	0.000000	0			No issue
gameresult	object	0.000000	0			No issue

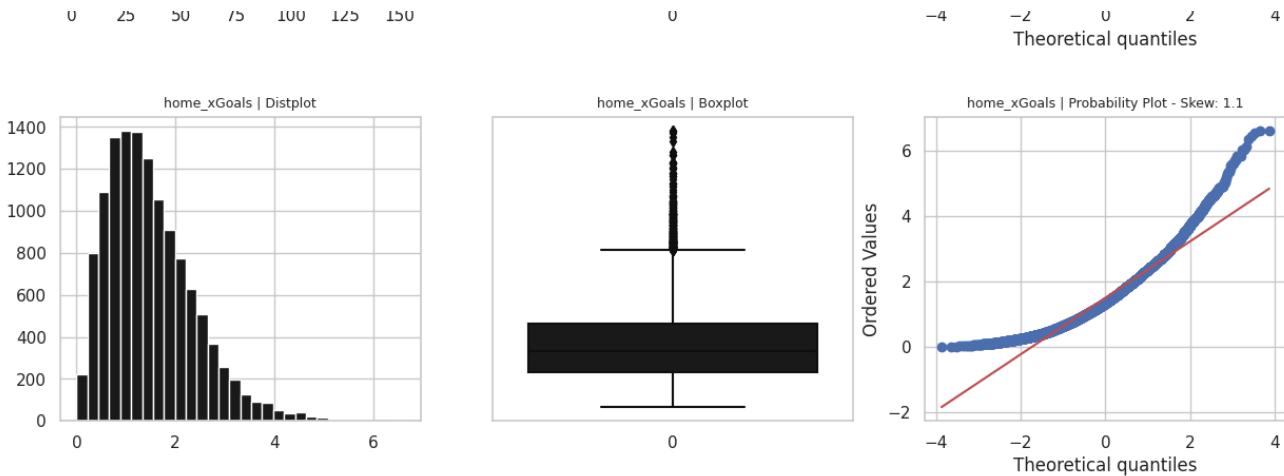
Number of All Scatter Plots = 21

Pair-wise Scatter Plot of all Continuous Variables

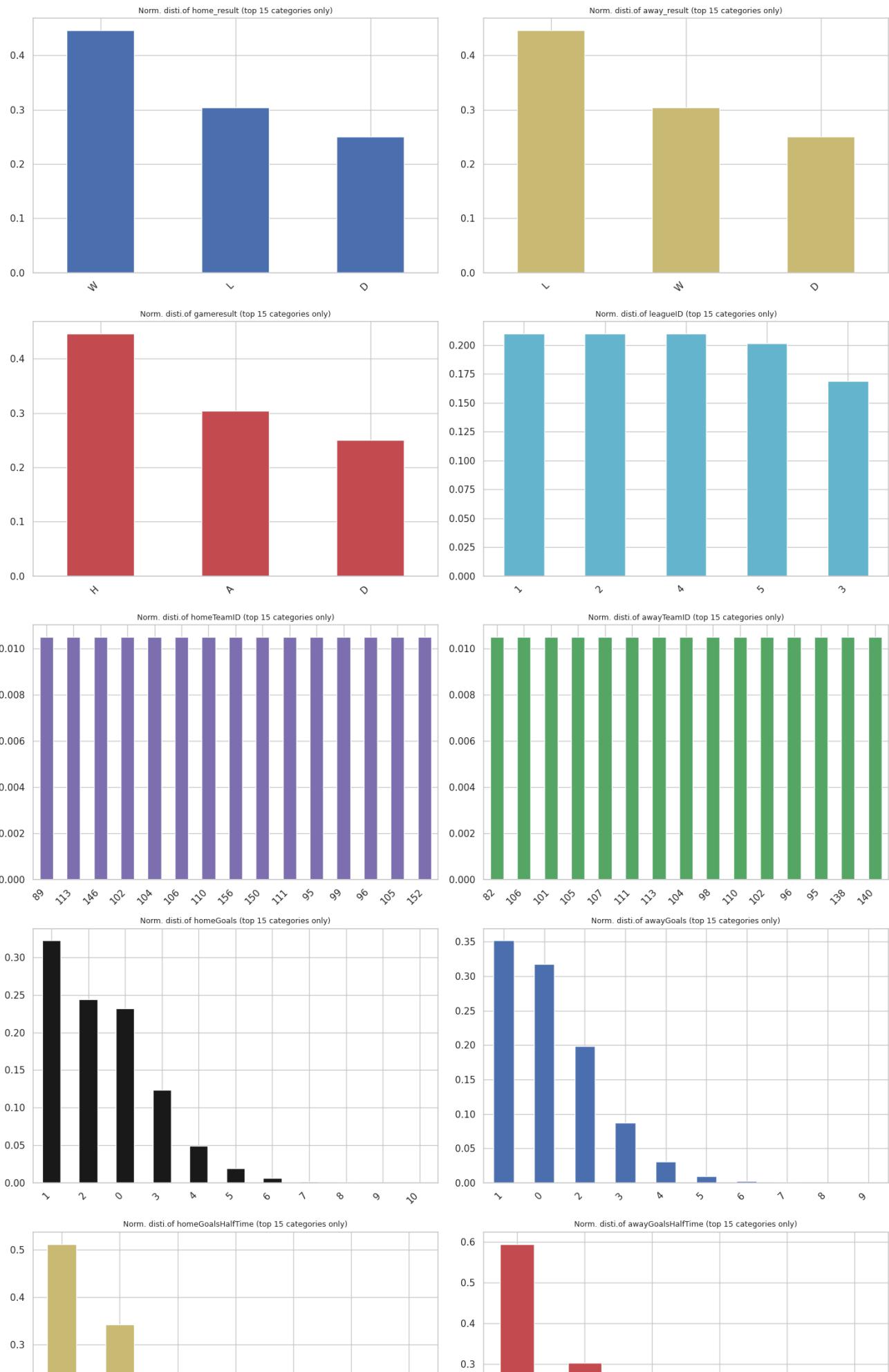


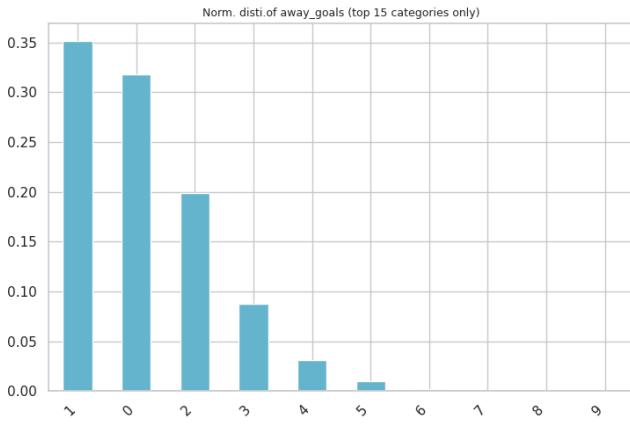
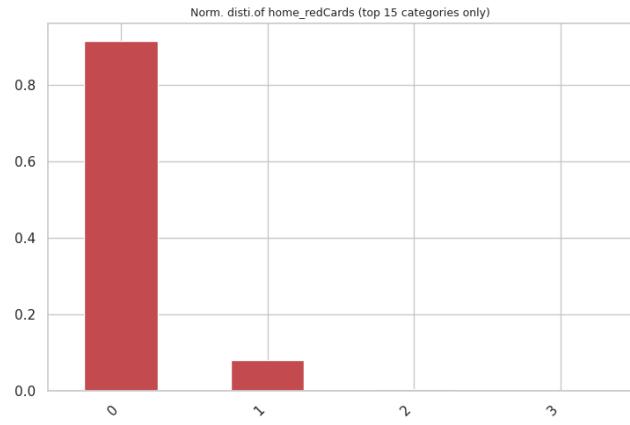
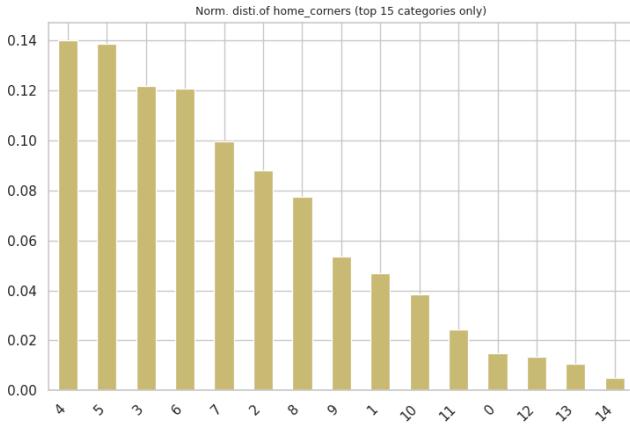
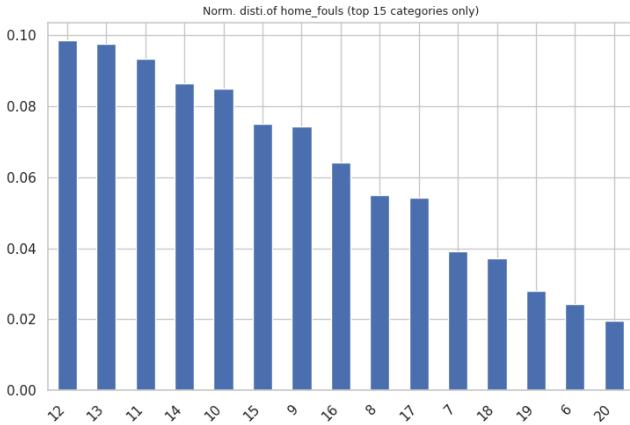
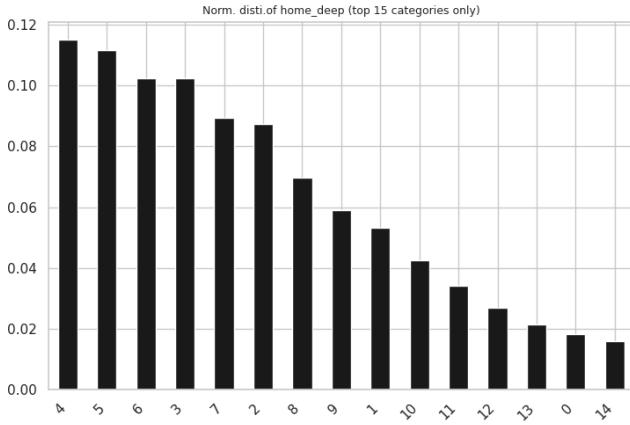
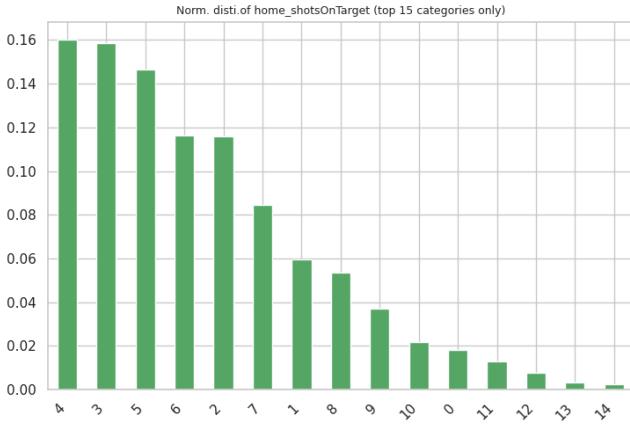
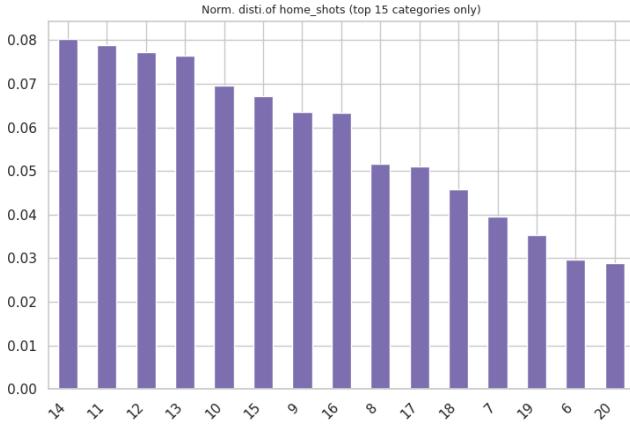
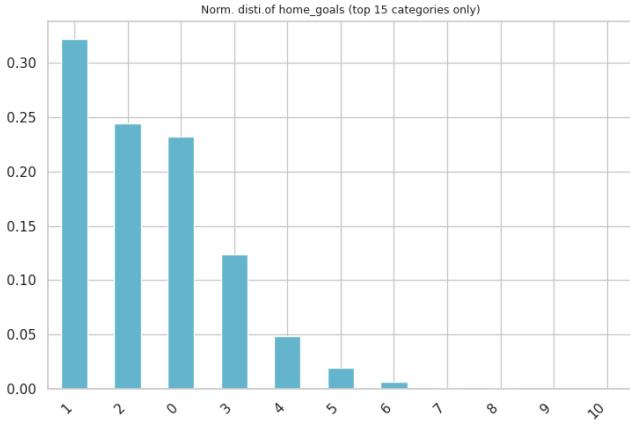
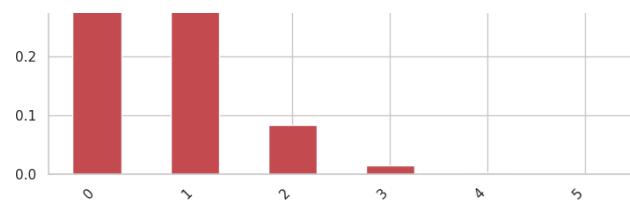
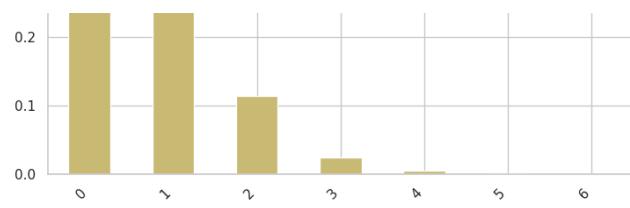




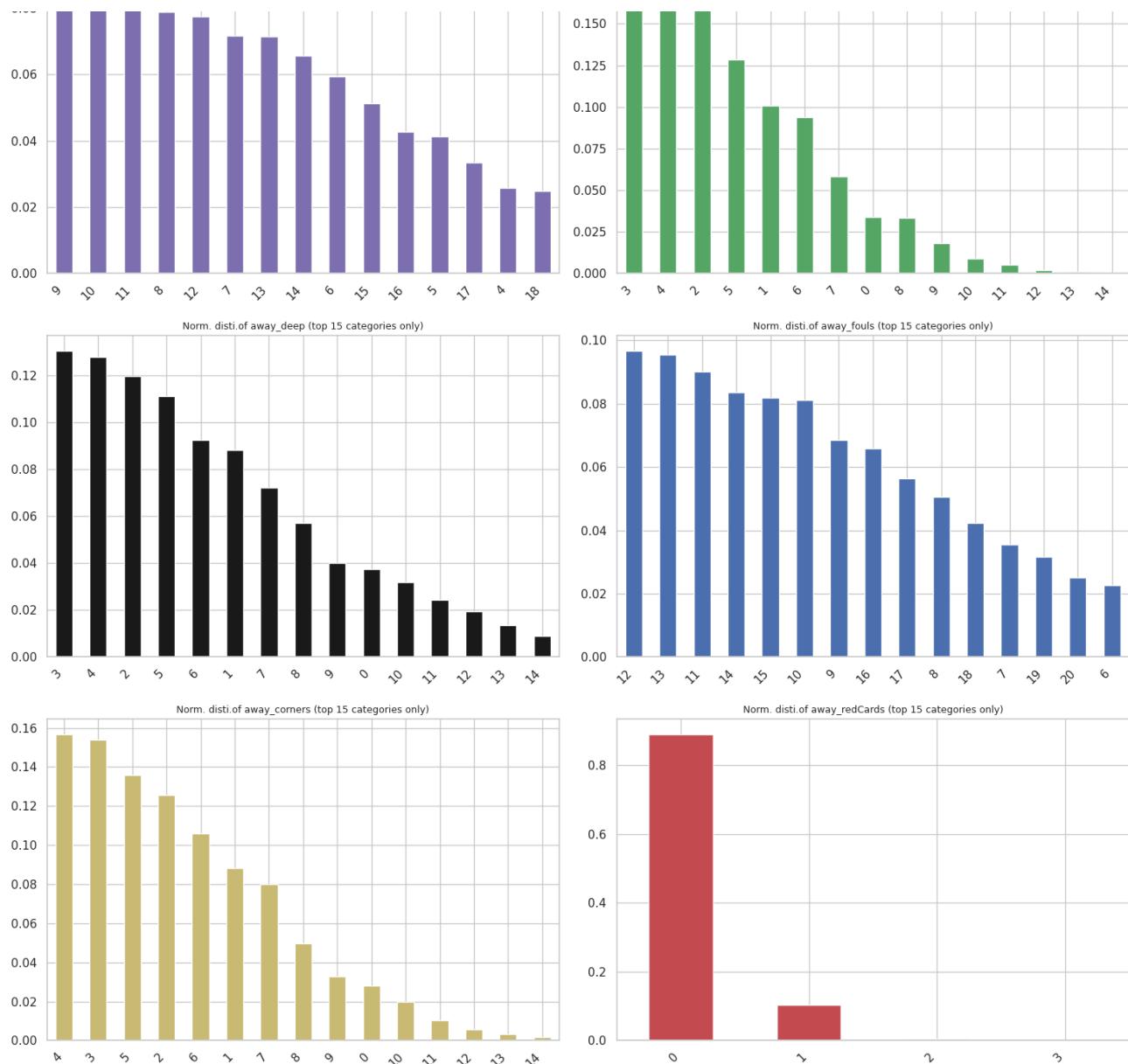


Histograms and Normalized distributions of all variables

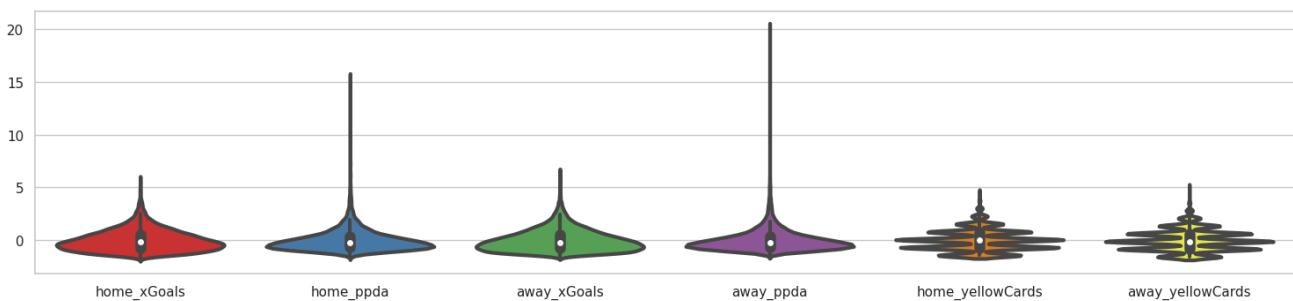


Notebook

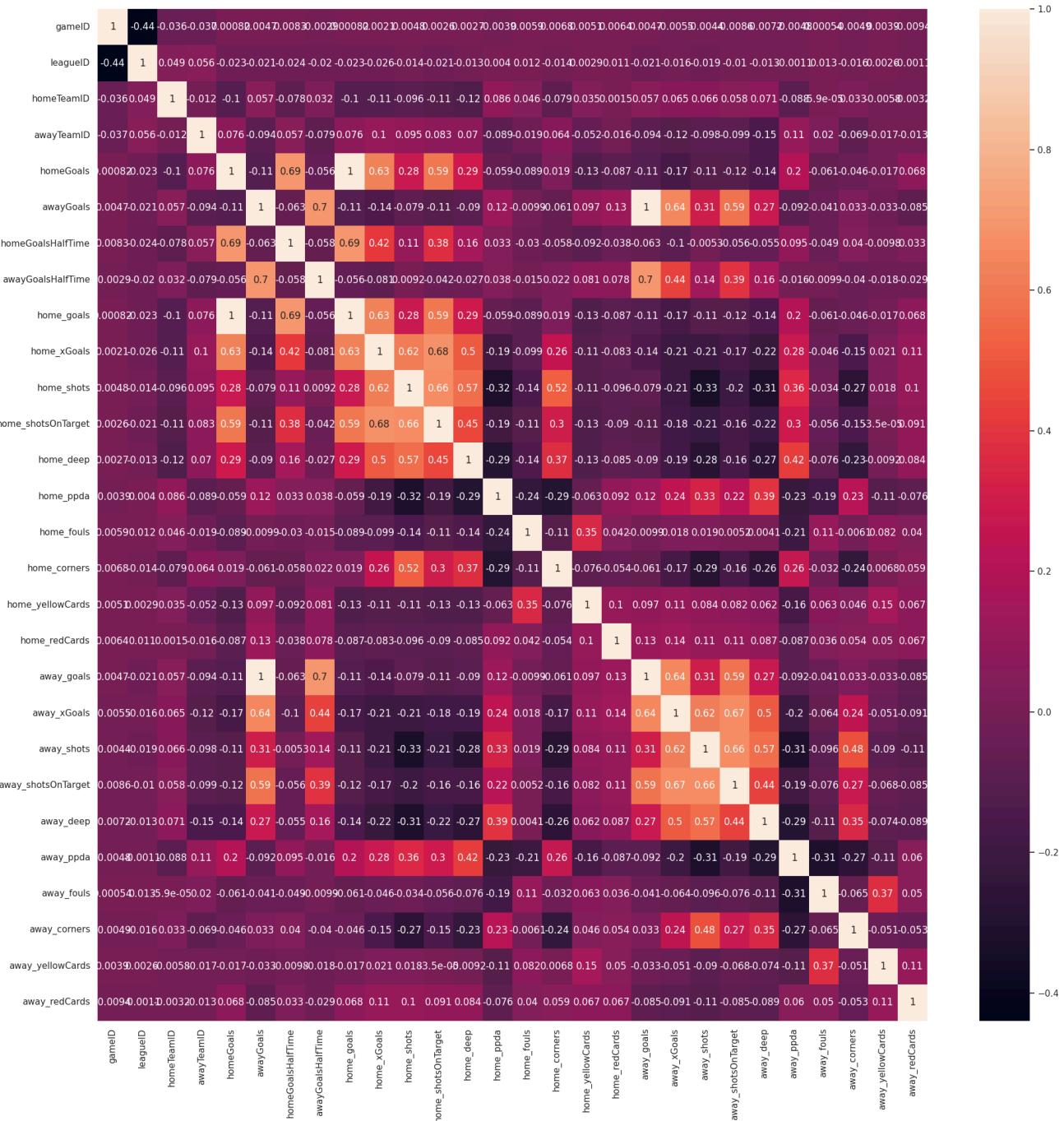
Notebook



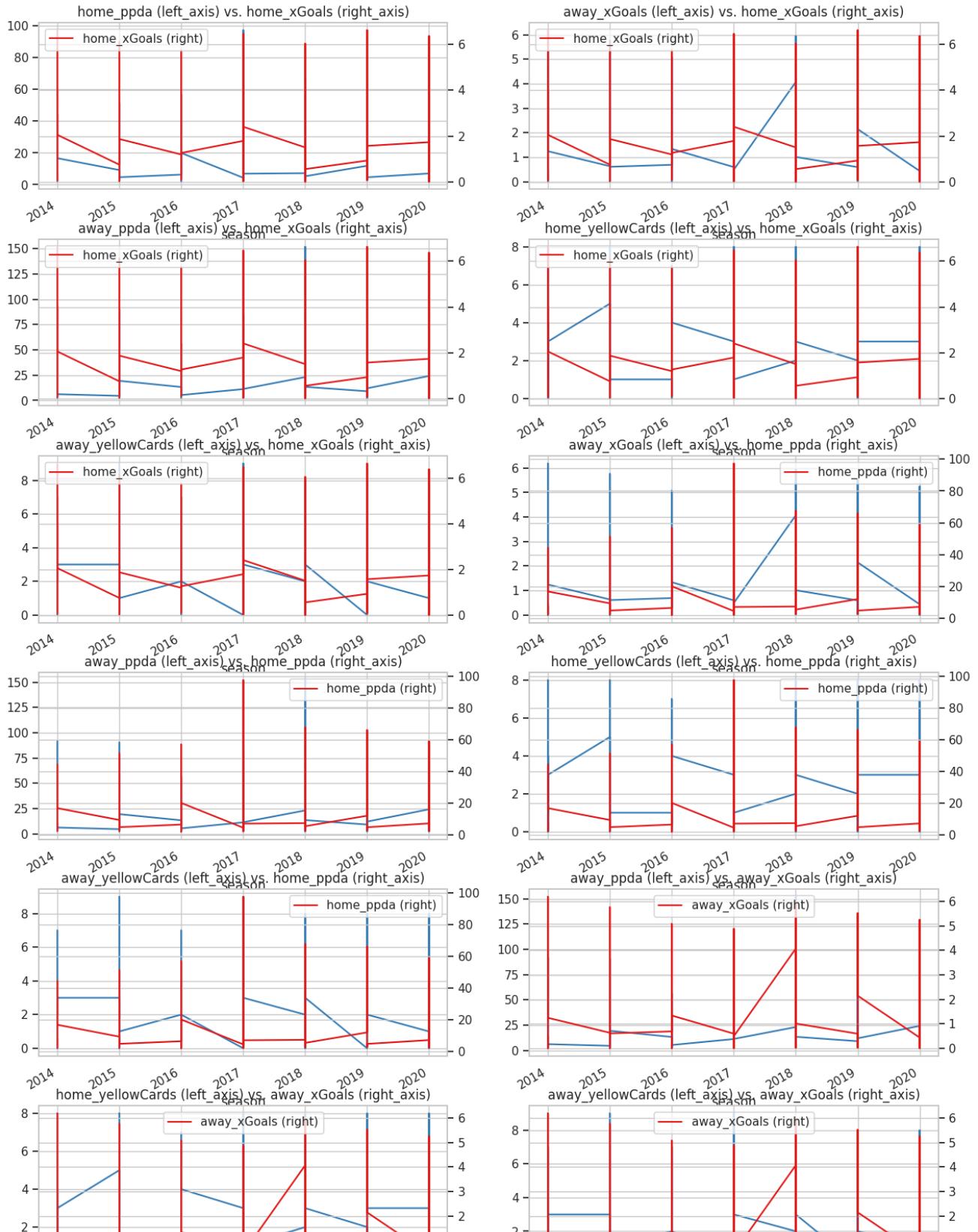
Violin Plot of all Continuous Variables



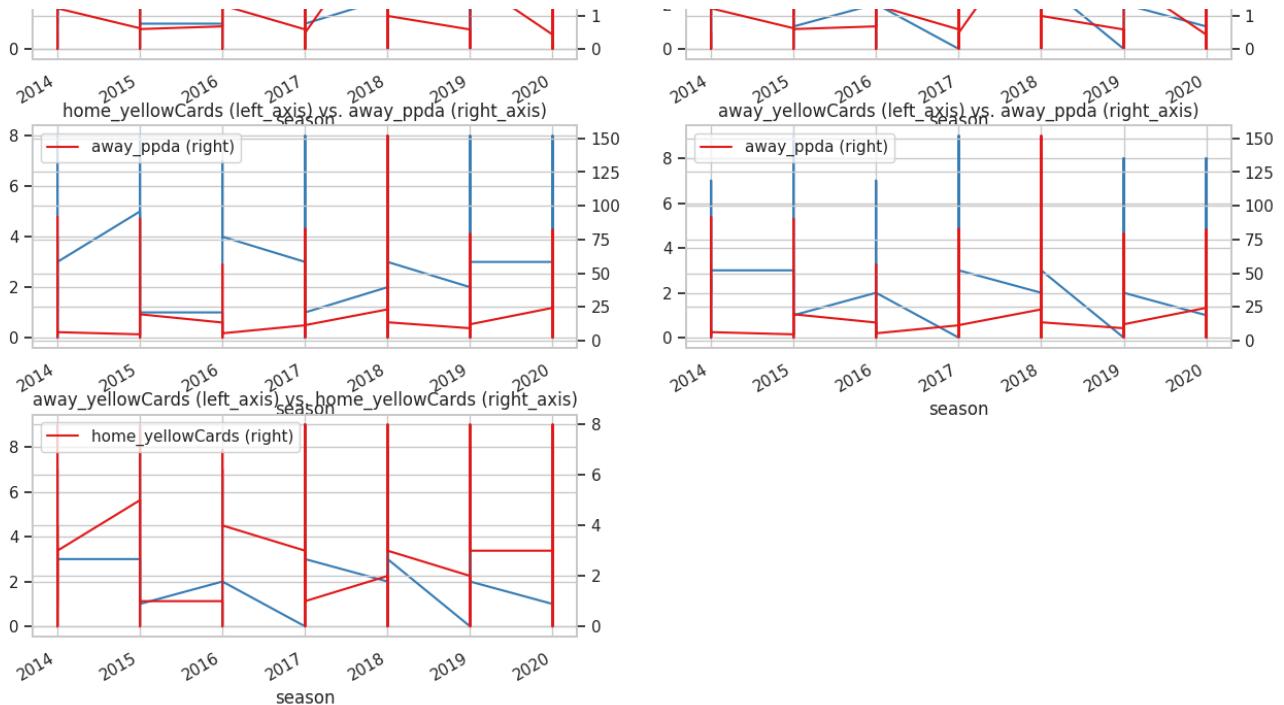
Time Series Data: Heatmap of Differenced Continuous vars including target =



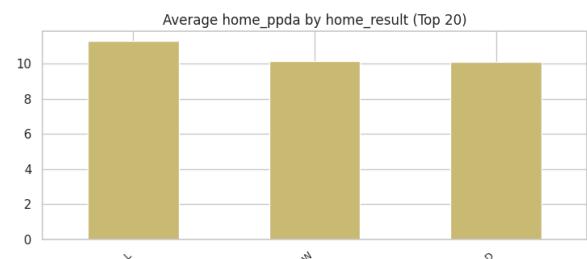
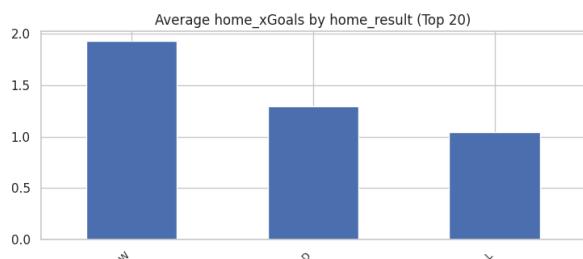
Time Series Plot by season: Pairwise Continuous Variables

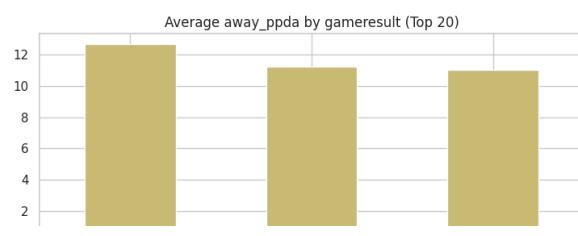
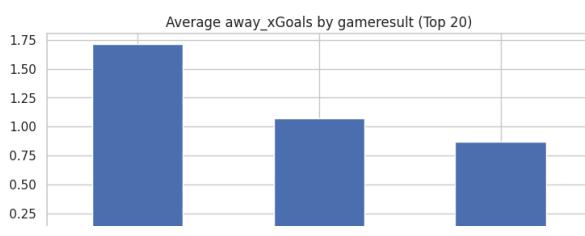
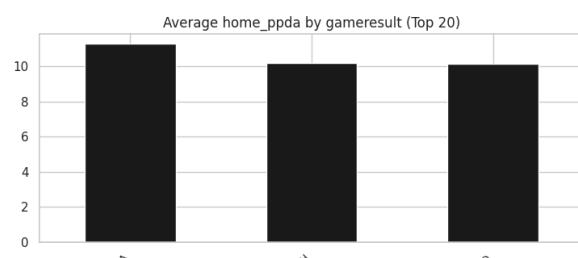
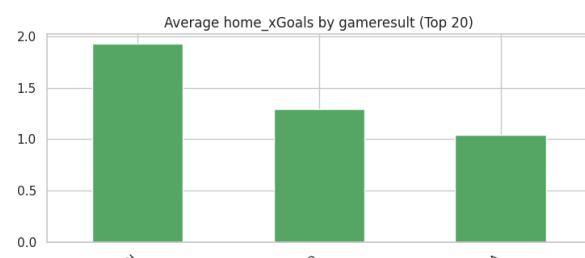
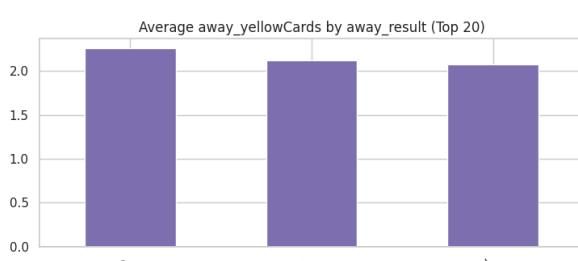
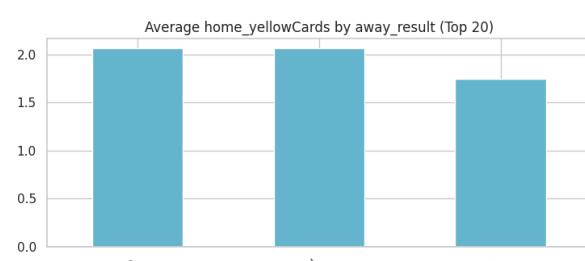
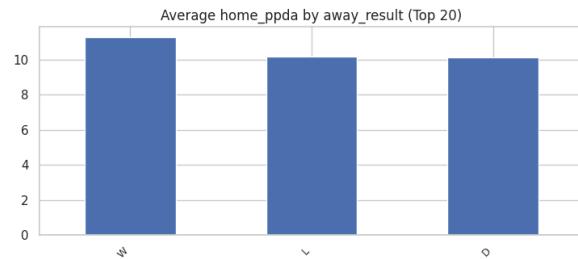
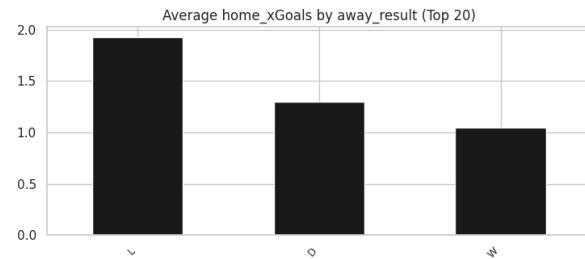
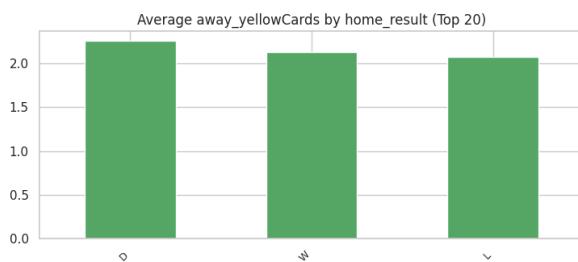
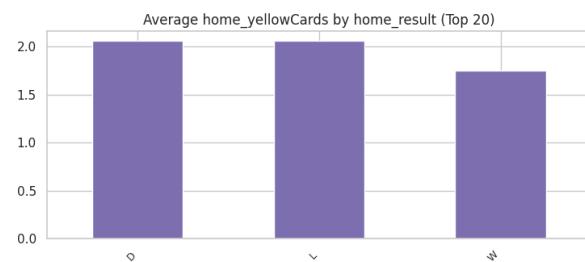
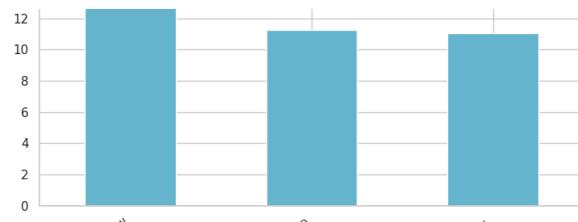
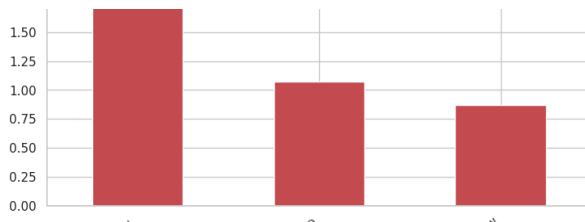


Notebook

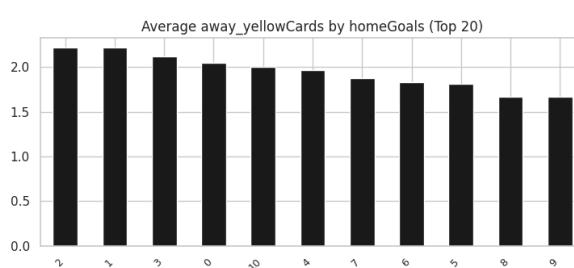
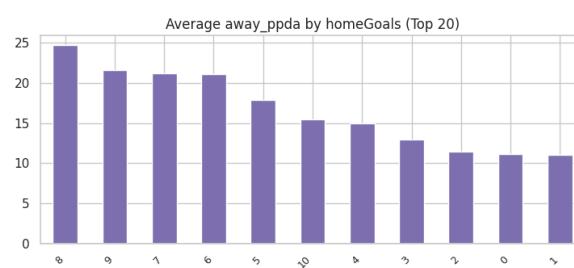
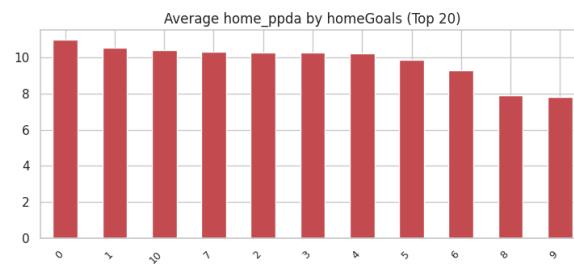
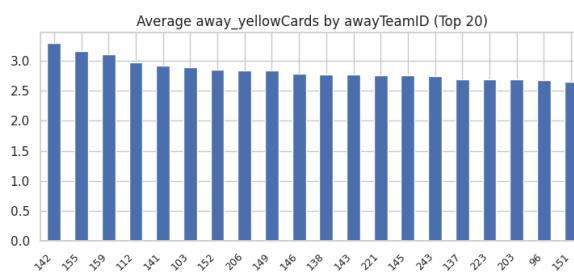
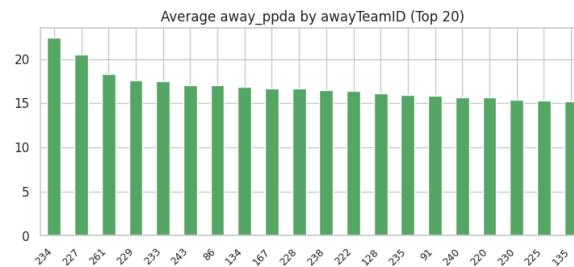
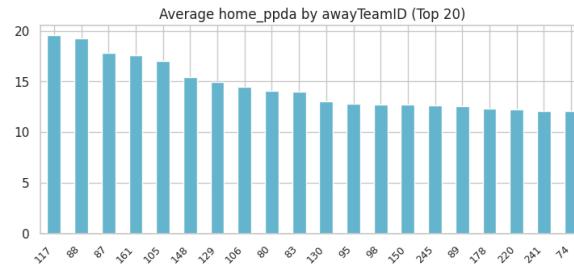
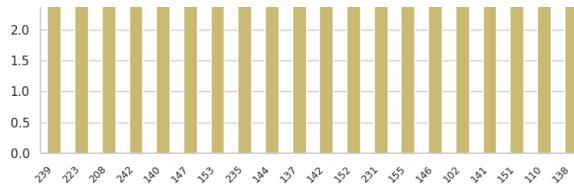
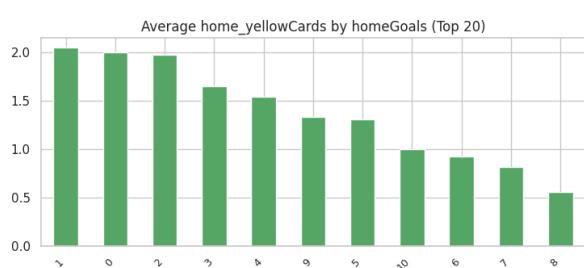
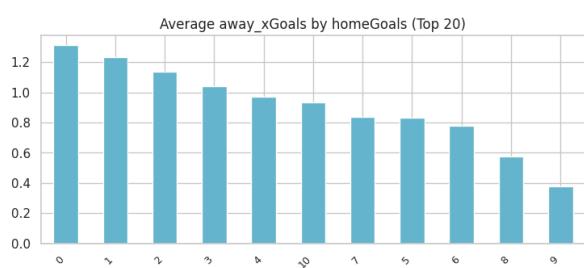
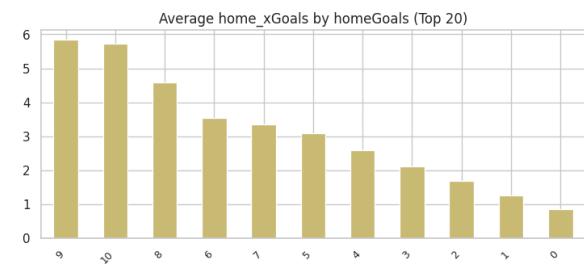
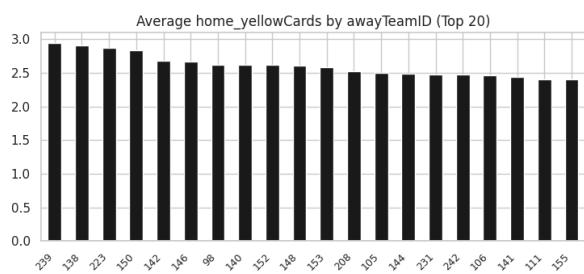
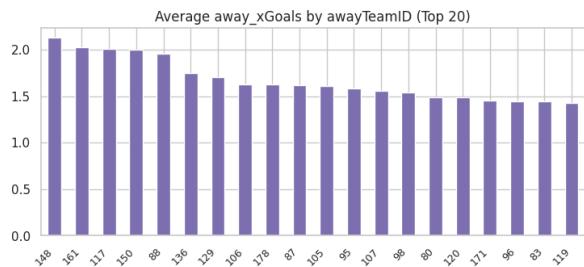
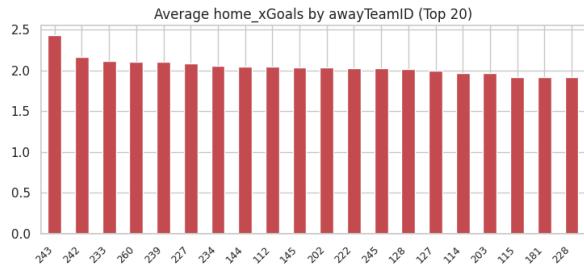
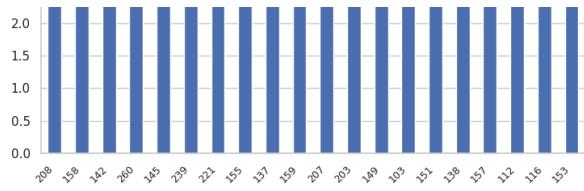


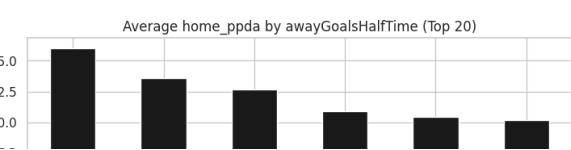
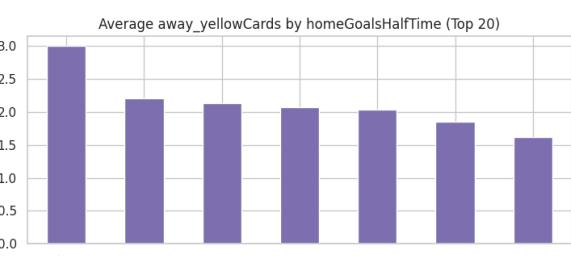
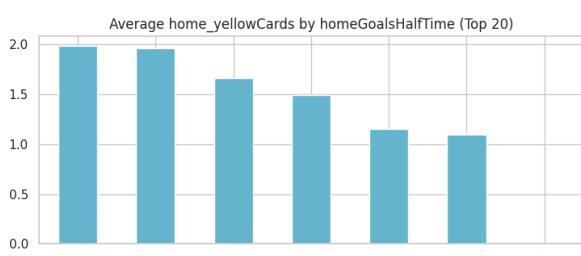
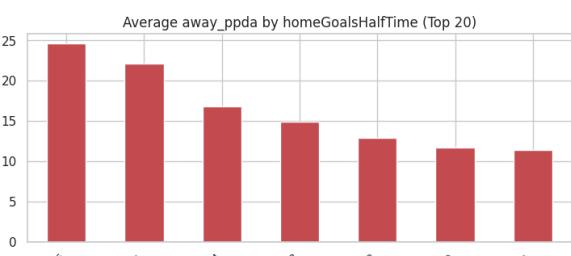
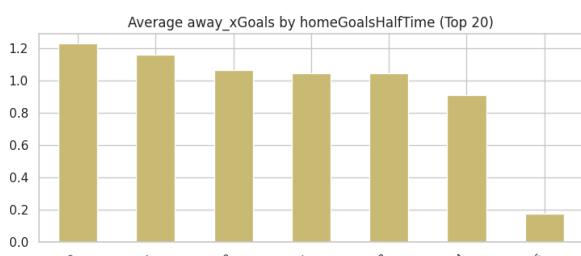
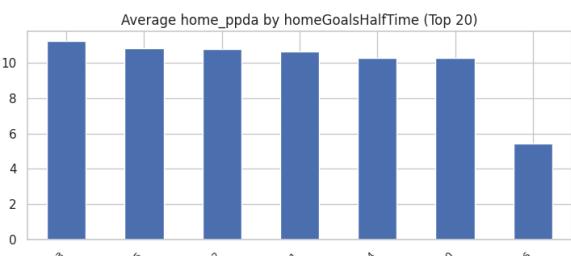
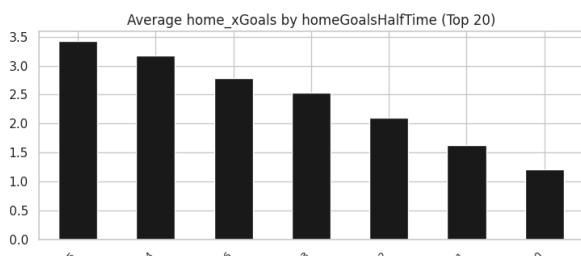
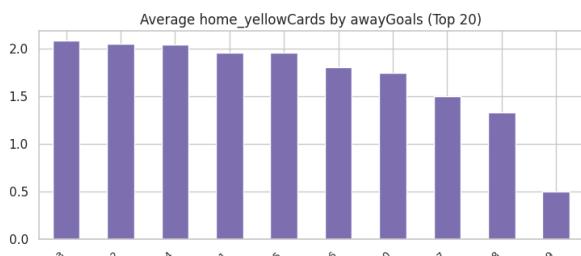
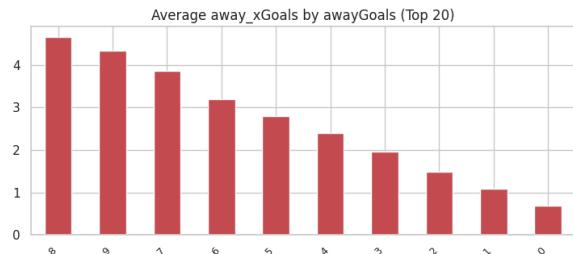
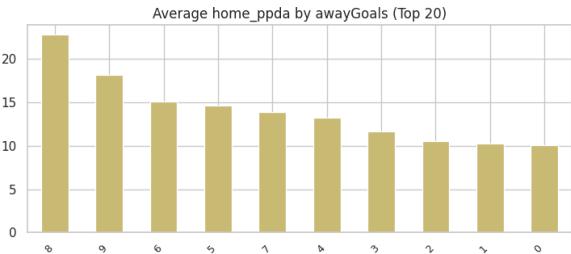
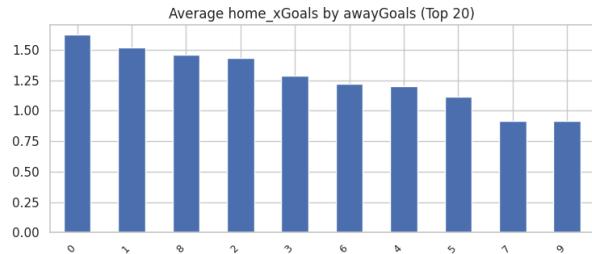
Bar plots for each Continuous by each Categorical variable

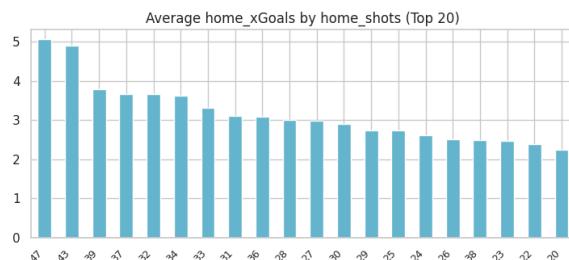
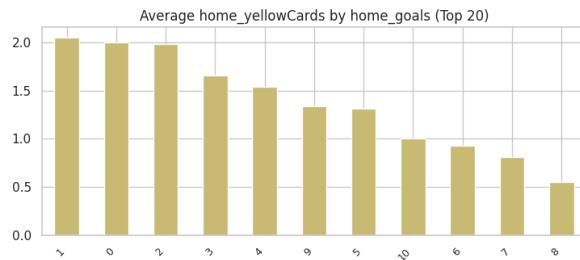
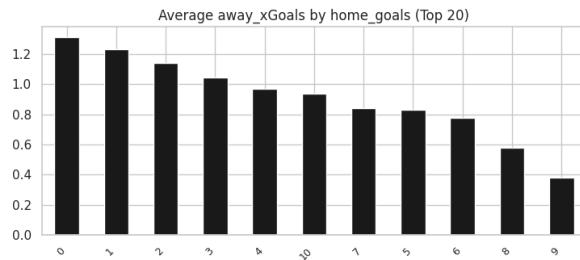
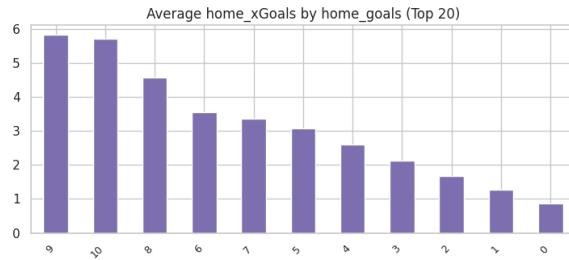
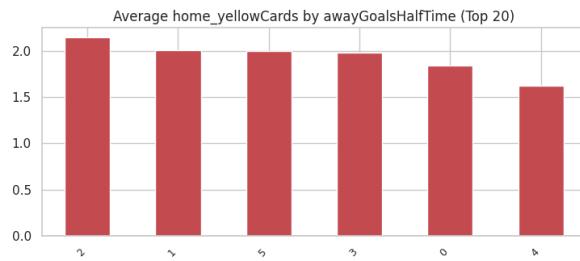


Notebook

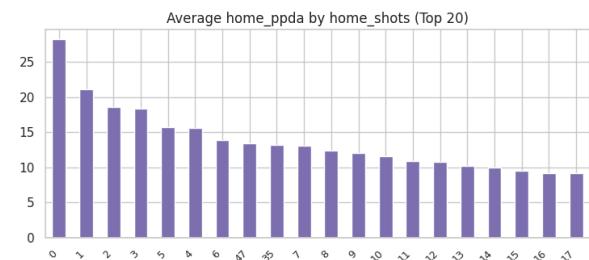
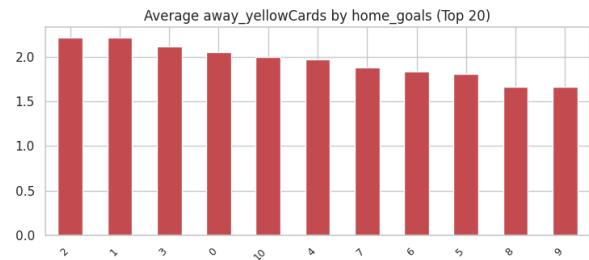
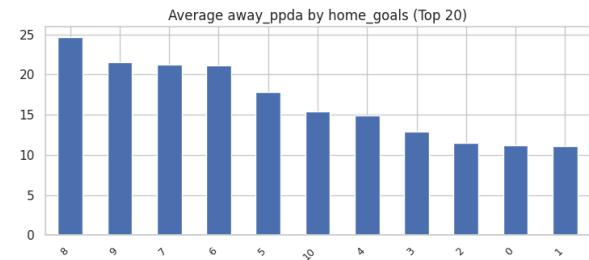
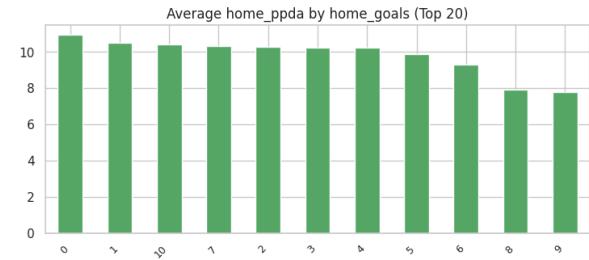
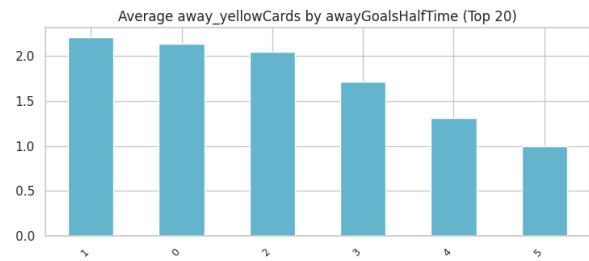
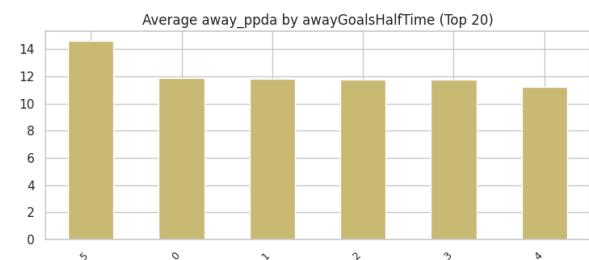
Notebook



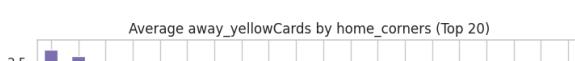
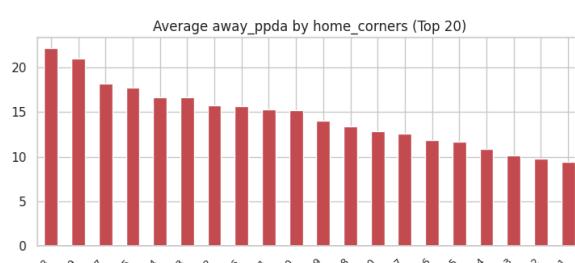
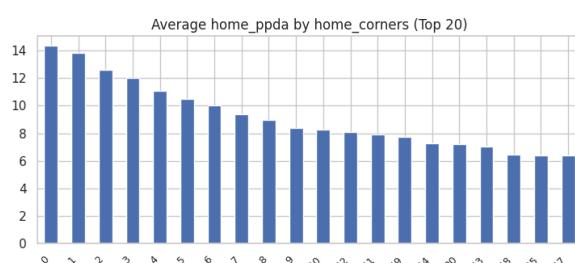
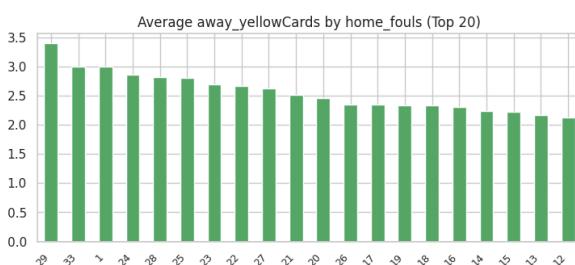
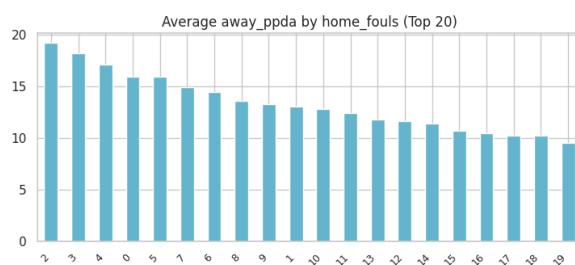
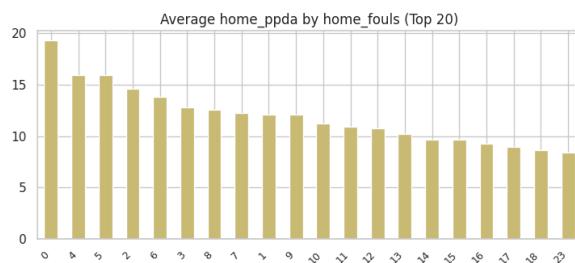
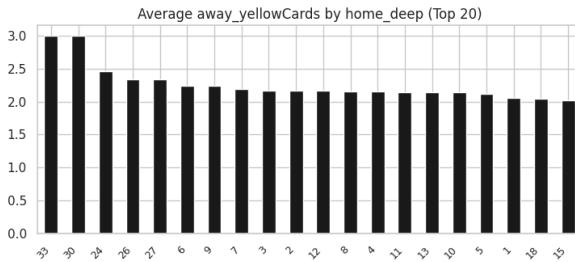
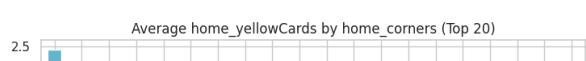
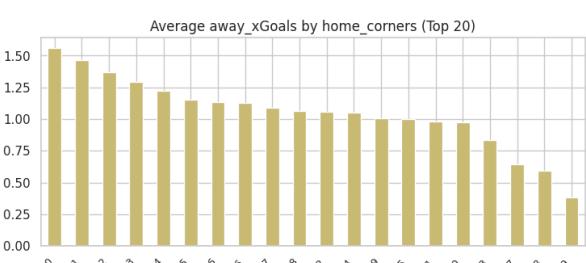
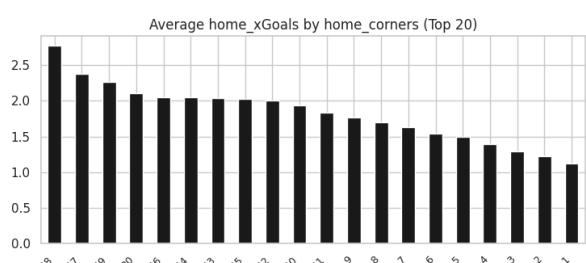
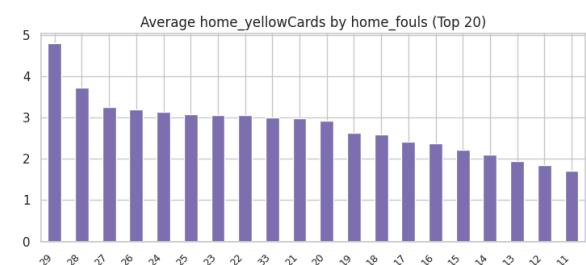
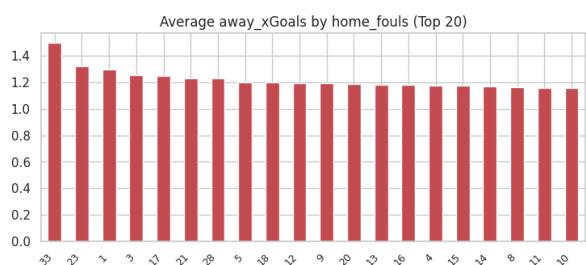
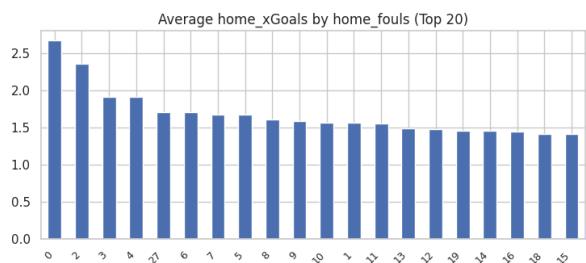
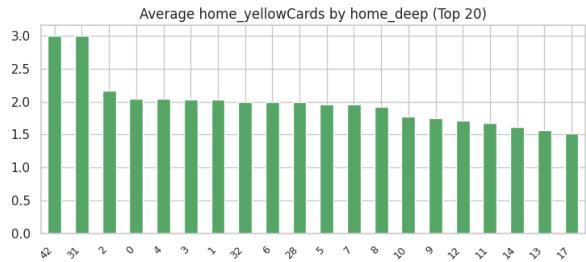
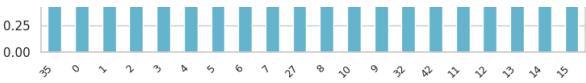


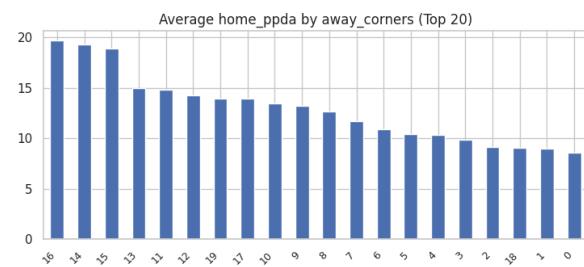
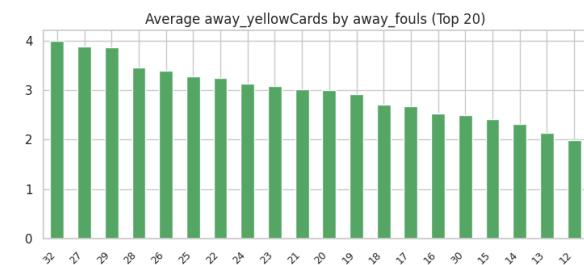
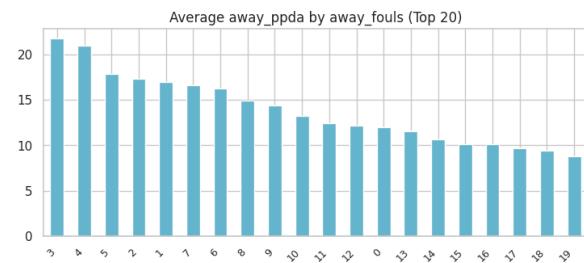
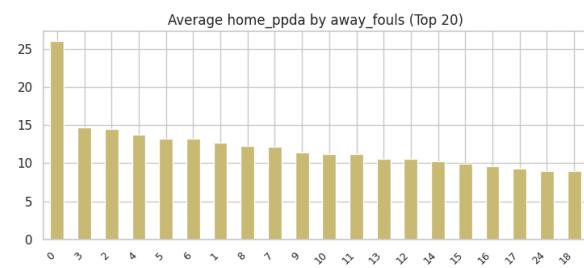
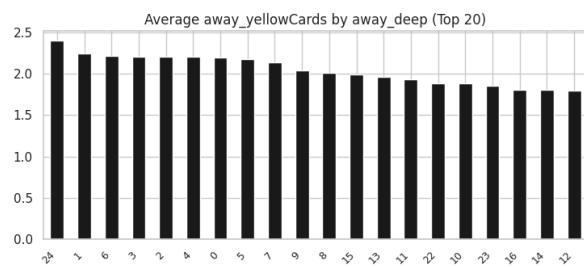
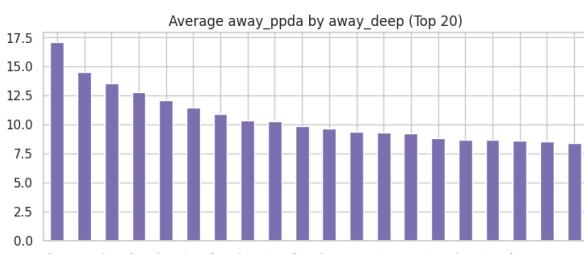
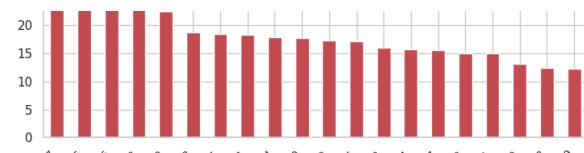
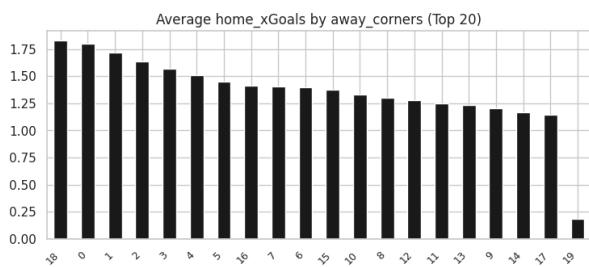
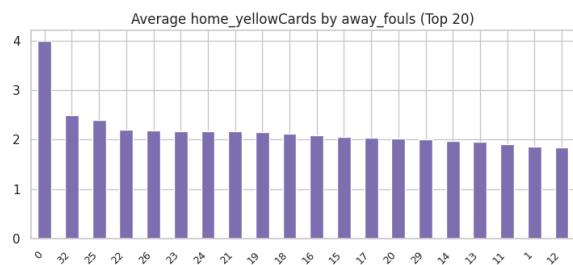
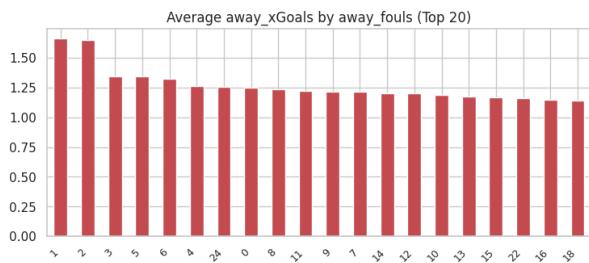
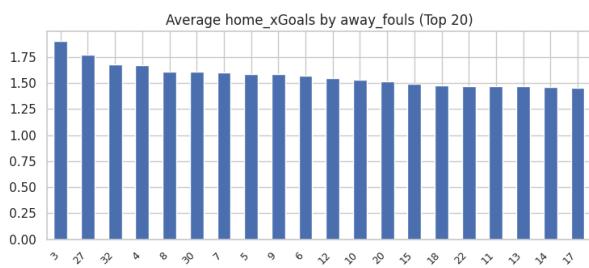
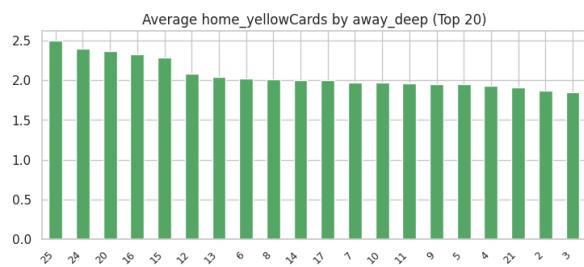
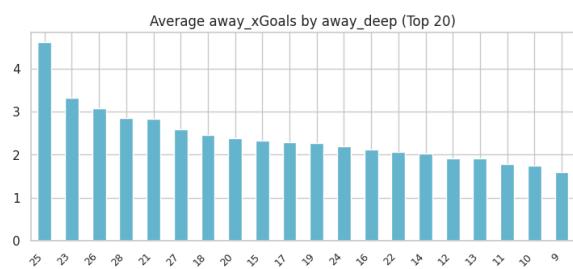
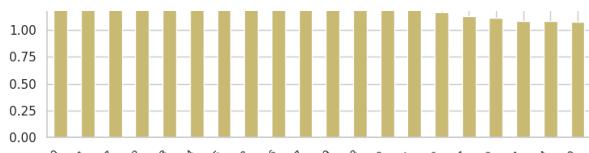


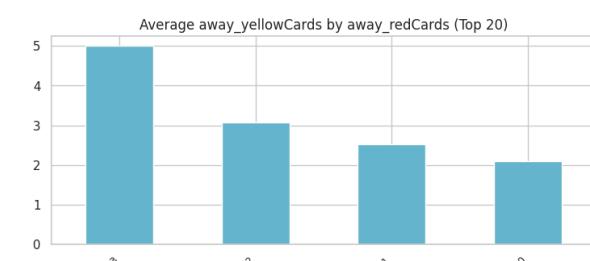
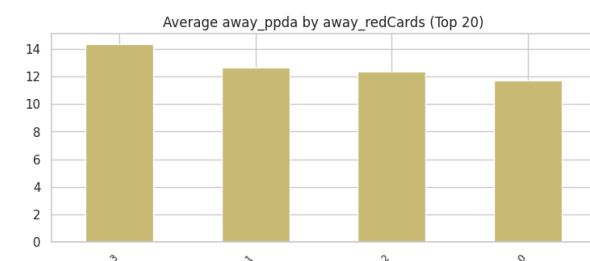
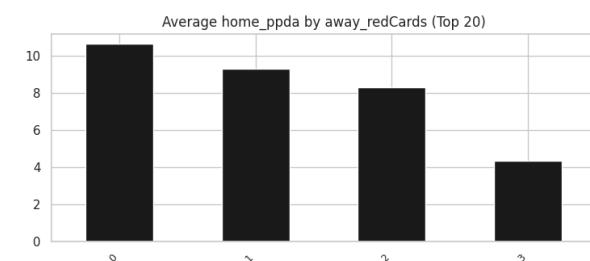
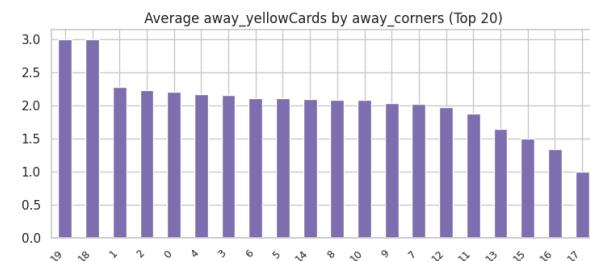
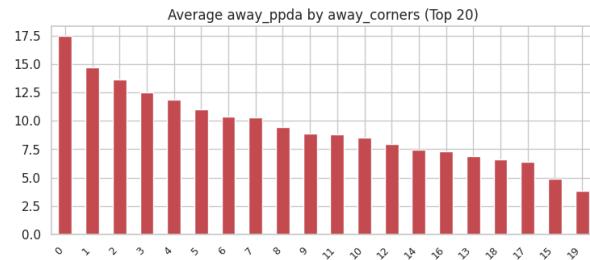
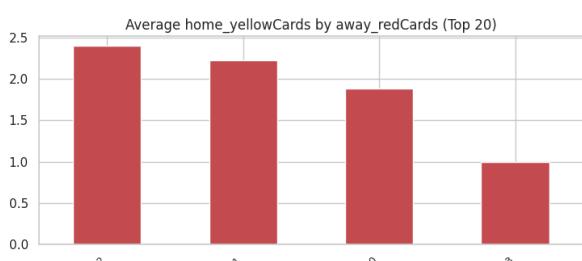
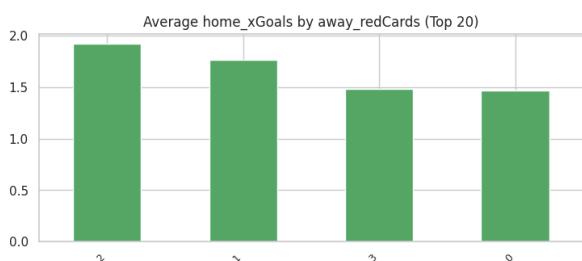
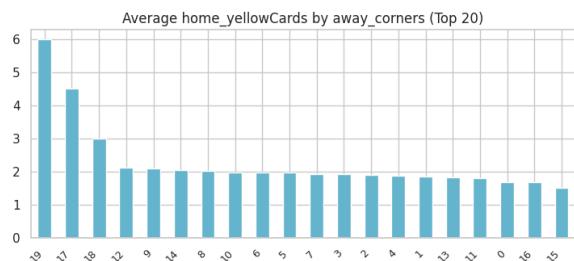
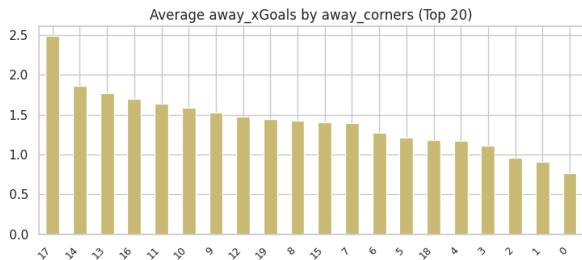
Notebook



Notebook



Notebook



```
[nltk_data] Downloading collection 'popular'
[nltk_data]   | Downloading package cmudict to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package cmudict is already up-to-date!
[nltk_data]   | Downloading package gazetteers to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package gazetteers is already up-to-date!
[nltk_data]   | Downloading package genesis to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package genesis is already up-to-date!
[nltk_data]   | Downloading package gutenberg to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package gutenberg is already up-to-date!
[nltk_data]   | Downloading package inaugural to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package inaugural is already up-to-date!
[nltk_data]   | Downloading package movie_reviews to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package movie_reviews is already up-to-date!
[nltk_data]   | Downloading package names to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package names is already up-to-date!
[nltk_data]   | Downloading package shakespeare to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package shakespeare is already up-to-date!
[nltk_data]   | Downloading package stopwords to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package stopwords is already up-to-date!
[nltk_data]   | Downloading package treebank to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package treebank is already up-to-date!
[nltk_data]   | Downloading package twitter_samples to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package twitter_samples is already up-to-date!
[nltk_data]   | Downloading package omw to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package omw is already up-to-date!
[nltk_data]   | Downloading package omw-1.4 to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package omw-1.4 is already up-to-date!
[nltk_data]   | Downloading package wordnet to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package wordnet is already up-to-date!
[nltk_data]   | Downloading package wordnet2021 to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package wordnet2021 is already up-to-date!
[nltk_data]   | Downloading package wordnet31 to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package wordnet31 is already up-to-date!
[nltk_data]   | Downloading package wordnet_ic to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package wordnet_ic is already up-to-date!
[nltk_data]   | Downloading package words to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package words is already up-to-date!
[nltk_data]   | Downloading package maxent_ne_chunker to
[nltk_data]   |     /home/leoadmin/nltk_data...
[nltk_data]   |     Package maxent_ne_chunker is already up-to-date!
[nltk_data]   | Downloading package punkt to
```

```
[nltk_data]      /home/leoadmin/nltk_data...
[nltk_data]      Package punkt is already up-to-date!
[nltk_data]      Downloading package snowball_data to
[nltk_data]          /home/leoadmin/nltk_data...
[nltk_data]      Package snowball_data is already up-to-date!
[nltk_data]      Downloading package averaged_perceptron_tagger to
[nltk_data]          /home/leoadmin/nltk_data...
[nltk_data]      Package averaged_perceptron_tagger is already up-
[nltk_data]          to-date!
[nltk_data]
[nltk_data] Done downloading collection popular
Could not draw wordcloud plot for date
Could not draw wordcloud plot for home_date
Could not draw wordcloud plot for away_date
All Plots done
Time to run AutoViz = 54 seconds

##### AUTO VISUALIZATION Completed #####
== Completed AutoViz on df_combined ==
```

